# SOME NEW RESULTS ON THE MAXIMUM GROWTH FACTOR IN GAUSSIAN ELIMINATION

ALAN EDELMAN AND JOHN URSCHEL

ABSTRACT. This paper combines modern numerical computation with theoretical results to improve our understanding of the growth factor problem for Gaussian elimination. On the computational side we obtain lower bounds for the maximum growth for complete pivoting for $n = 1 : 75$ and $n = 100$ using the Julia JuMP optimization package. At $n = 100$ we obtain a growth factor bigger than $3n$. The numerical evidence suggests that the maximum growth factor is bigger than $n$ if and only if $n \geq 11$.

We also present a number of theoretical results. We show that the maximum growth factor over matrices with entries restricted to a subset of the reals is nearly equal to the maximum growth factor over all real matrices. We also show that the growth factors under floating point arithmetic and exact arithmetic are nearly identical. Finally, through numerical search, and stability and extrapolation results, we provide improved lower bounds for the maximum growth factor. Specifically, we find that the largest growth factor is bigger than $1.0045n$ for $n > 10$, and the lim sup of the ratio with $n$ is greater than or equal to 3.317. In contrast to the old conjecture that growth might never be bigger than $n$, it seems likely that the maximum growth divided by $n$ goes to infinity as $n \to \infty$.

## 1. INTRODUCTION

We begin with a sketch of the history of the subject. For an introduction to the technical background of the field, and a more technical discussion of related work, see Subsections 1.3 and 1.4, respectively.

### 1.1. **History of Complete Pivoting (Overview).** 

Understanding and bounding the growth factor for Gaussian elimination has intrigued mathematical numerical linear algebraists for many decades. It is one of those beautiful problems that is quite easy to state, and yet we still know so very little, especially in the context of complete pivoting. This may seem even more surprising as the problem has been around for more than six decades.

**The 1960s-1970s:** In 1961, Wilkinson [35, Equation 4.15] published a famous bound for the growth factor for complete pivoting which was always considered pessimistic (see Eq. (1.1) of this paper). In that same paper he writes "no matrix has been encountered for which [the growth factor] was as large as 8."

In 1964, Leonard Tornheim wrote in a technical paper [28] "there is a conjecture, attributed to J.H. Wilkinson, that" the growth factor for complete pivoting would be at most $n$. This is the first written record of this conjecture (known to be false since 1991) as far as we are aware. Leonard Tornheim's work is notable in that it was not published in any journals, but rather as technical papers and Chevron research reports during a time when Chevron had a major presence in Richmond, CA, some 20 miles from UC Berkeley.[1] One year later, in 1965, Wilkinson published the famous book *The Algebraic Eigenvalue Problem* [36, p.213] in which he wrote that no matrix had yet been discovered with growth greater than $n$ for an $n \times n$ matrix. In that same year a cover letter written by Eugene B. Reid[2] for a report by Tornheim [29] writes (without any reference) that it was a *widely known conjecture* that growth would always be less than or equal to $n$. In 1968, Cryer [4] specifically references Wilkinson's 1965 words as a conjecture, though Wilkinson never explicitly put to print a formal conjecture. Nonetheless Cryer specifically wrote in his abstract "It has been conjectured by Wilkinson...". A possible best guess is that through the rumor grape vine often known as the children's game of "broken telephone," the observation that none was ever seen morphed into a folk conjecture and thus we will attribute this conjecture not to Wilkinson (who did not write it down) nor Cryer (who wrote it later), but to folklore.[3]

The 1960s saw the maximum growth computed for $n = 1, 2, 3, 4$ and bounded for $n = 5$ by the growth chasers Tornheim [28, 29, 30, 31], Cryer [4], and Cohen [2]. Hadamard matrices (matrices with all entries $\pm 1$ and orthogonal columns) were shown by Tornheim [28] and Cryer [4] to have last pivot $n$. It became natural to wonder whether Hadamard matrices could be a counterexample to the Wilkinson conjecture.

**The 1980s:** After a bit of a lull, in 1988 Day and Peterson wrote a lovely article in the American Math Monthly [5] which revisits $n = 3$, studies Hadamard matrices, and notably is the first to explore the growth problem for complete pivoting with numerical optimization software, specifically the NPSOL Library out of Stanford (Nonlinear Programming, Stanford Optimization Laboratory). In particular they were the first to observe the number 4.1325 when $n = 5$ as an output of the optimization software.

---

[1]Tornheim was a rather active mathematician. In a March 2023 phone call, the first author contacted Tornheim's son who described how Tornheim wrote a computer program that made Chevron ten million dollars, thus establishing in his words "what a mathematician could be good for" (in private industry). He also mentioned that his father had moved to the east coast and as coincidence would have it, he lived directly across the street in Brookline, MA but sadly, the first author did not know this and recognizes what a lost opportunity this might have been.

[2]In March 2023, the first author telephoned the son of Eugene B. Reid who described his father as having bought the first commercial computer on the West Coast of the US, and that his father was not a mathematician or a statistician, but he was responsible for the computer, mathematical, and statistical activities at Chevron. He also bought Chevron's very first computer which his son claims might have been the first commercial computer on the West Coast of the US.

[3]In another March 2023 phone call, the first author spoke to William Kahan who seemed rather certain that Wilkinson had never stated the conjecture in so many words. Kahan described how computers were slow enough at the time that Wilkinson would watch the bits go by and look for large growth, and growth $> n$ was never observed. Nonetheless, we can not completely rule out that Wilkinson may have verbally stated or at least hinted at the conjecture. In fact, Cleve Moler stated publicly on August 16, 2023 in Oxford, UK that Wilkinson was not inclined to use terms such as "theorem," "proof," and "conjecture," but Moler felt that Wilkinson had believed the conjecture.

In 1989, Higham and Higham [18] pointed out that many common matrices can have growth factors of order $n$ (for any pivoting strategy).

**The early 1990s:** Interest in the growth factor was substantially rejuvenated when Trefethen and Schreiber [34] performed average case analyses of the growth factor in 1990. One year later, Nick Gould [15] surprised everyone by finding a 13x13 matrix with growth bigger than 13 in finite precision using his LANCELOT software. The solution was confirmed to be near a true example in exact arithmetic in 1992 [7].

**1993-Present:** In the over 30 years since, there was no progress whatsoever in improving Gould's numbers for complete pivoting through computation (which would raise a lower bound) or lowering any mathematical upper bounds. This is a testament to the difficulty of the problem.

## 1.2. **Other pivoting analyses.**

**No Pivoting:** In 2006, the celebrated smoothed analysis[4] of Sankar, Spielman, and Teng [26] showed that large growth is unlikely from a probabilistic perturbative viewpoint with no pivoting, and pointed out that such an analysis could be possible for partial and complete pivoting.

**Partial Pivoting:** In 1994, Foster [11] pointed out that practical problems can bump into the unacceptable $2^{n-1}$ bound for partial pivoting. The first author remarked in [9] that numerical experiments suggested in contrast to [34] that the growth might be more like $O(n^{1/2})$ than $O(n^{2/3})$ on average. In addition to the smoothed analysis for no pivoting, Sankar [25] also performed a smoothed analysis of partial pivoting with subexponential bounds. Very recently Huang and Tikhomirov [19] obtained new results exploring the average case analysis for partial pivoting.

**Complete Pivoting for Hadamard Matrices:** It remains unknown, though perhaps it seems unlikely, that a Hadamard matrix could have an earlier pivot bigger than $n$, given that the last three pivots can only be $n/2, n/2$ and $n$, and the fourth from the end is at most $n/2$. Nonetheless, complete pivot patterns for Hadamard matrices remain a fascinating topic of research. A comprehensive review of the topic including new progress written in 2013 by Kravvaritis may be found in [21]. Of note are the investigations by Seberry [27] and also [8, 9]. We note that the growth factor for a Hadamard matrix of dimension $n \leq 16$ is known to be $n$ under complete pivoting.

1.3. **Technical Background.** The solution of a linear system, i.e., given a matrix $A$ and vector $b$, finding a vector $x$ satisfying $Ax = b$, is one of the oldest problems in mathematics. Gaussian elimination, a technique in which a matrix is factored into the product of a lower and upper triangular matrix, is one of the most fundamental and important techniques for solving linear systems. The algorithm proceeds by converting

---

[4]Incidentally "smoothed analysis" was named by the first author in his car while driving Dan Spielman and Shanghua Teng in Cambridge, MA.

$A$ into upper triangular form through row operations. In particular, given an $n \times n$ matrix $A = (a_{i,j})$, Gaussian elimination performs the iteration

$$a_{i,j}^{(1)} := a_{i,j} \qquad \text{for} \quad i,j = 1, ..., n,$$

$$a_{i,j}^{(k+1)} := a_{i,j}^{(k)} - \frac{a_{i,k}^{(k)} a_{k,j}^{(k)}}{a_{k,k}^{(k)}} \qquad \text{for} \quad i,j = k, ..., n, \ k = 1, ..., n-1.$$

This can be equivalently written as successive rank one updates of sub-matrices of $A$, i.e.,

$$A^{(k+1)} := A_{k+1:n,k+1:n}^{(k)} - \frac{1}{a_{k,k}^{(k)}} A_{k+1:n,k}^{(k)} A_{k,k+1:n}^{(k)} \qquad \text{for} \quad k = 1, ..., n-1,$$

where $A^{(k)} = (a_{i,j}^{(k)})_{i,j \geq k}$ and $A_{i_1:i_2,j_1:j_2}$ is defined as the sub-matrix of $A$ containing only rows $\{i_1, ..., i_2\}$ and columns $\{j_1, ..., j_2\}$. The resulting LU factorization of $A$ is given by

$$L(i,j) = \frac{a_{i,j}^{(j)}}{a_{j,j}^{(j)}} \quad \text{for} \quad i \geq j, \quad \text{and} \quad U(i,j) = a_{i,j}^{(i)} \quad \text{for} \quad j \geq i,$$

and this factorization is unique (up to scaling, i.e., $A = (LD)(D^{-1}U)$ for any invertible diagonal matrix $D$). Not all matrices have an LU factorization (issues arise if $a_{k,k}^{(k)} = 0$ for some $k < n$), and may require a permutation of the rows (or, equivalently, columns) of the matrix in order for such a factorization to exist. In addition, when computations are performed in finite precision, issues due to round-off error can occur. The backward error due to rounding in Gaussian elimination can be estimated by the number of bits of precision, the condition number of the matrix $A$, and the growth factor of the Gaussian elimination algorithm (see [20, Theorem 2.6] or [17, Theorem 9.5] for details). For this reason, understanding the growth factor under different permutation strategies is of both theoretical and practical importance. Using exact arithmetic, the growth factor of Gaussian elimination is defined as

$$g(A) := \frac{\max_{i,j,k} |a_{i,j}^{(k)}|}{\max_{i,j} |a_{i,j}|}.$$

When performing Gaussian elimination in finite precision, say, using only numbers that can be represented in base $\beta$ with a length $t$ mantissa, the algorithm suffers from round-off error, and the growth factor in this setting may be larger than $g(A)$. However, as we will see in Section 4, when $t = \omega(\log_\beta^2 n)$, the maximum growth factors in exact and floating point arithmetic are nearly identical (up to a $1 - o(1)$ multiplicative factor) under complete pivoting (see Theorem 4.2). For this reason, we focus almost exclusively (save for Section 4) on exact arithmetic. The most popular and well-studied methods for permuting a matrix in Gaussian elimination are partial pivoting (requiring $|a_{i,k}^{(k)}| \leq |a_{k,k}^{(k)}|$), complete pivoting (requiring $|a_{i,j}^{(k)}| \leq |a_{k,k}^{(k)}|$), and the slightly less well-known rook pivoting (requiring $|a_{i,k}^{(k)}|, |a_{k,j}^{(k)}| \leq |a_{k,k}^{(k)}|$). The growth factor for partial pivoting is well understood in the worst case, and so, in this work, we primarily focus on complete pivoting and, to some extent, rook pivoting as well.

Let $\mathbf{GL}_n(\mathbb{C})$ be the set of $n \times n$ non-singular complex matrices. For simplicity, when considering a given pivoting strategy, we simply restrict ourselves to the set of matrices

that satisfy the constraints of the pivoting procedure without requiring pivoting. In particular, we define

$$\mathbf{PP}_n(S) = \{A \in \mathbf{GL}_n(\mathbb{C}) \cap S^{n \times n} \,|\, |a_{i,k}^{(k)}| \le |a_{k,k}^{(k)}| \text{ for all } i \ge k\},$$

$$\mathbf{CP}_n(S) = \{A \in \mathbf{GL}_n(\mathbb{C}) \cap S^{n \times n} \,|\, |a_{i,j}^{(k)}| \le |a_{k,k}^{(k)}| \text{ for all } i, j \ge k\},$$

$$\mathbf{RP}_n(S) = \{A \in \mathbf{GL}_n(\mathbb{C}) \cap S^{n \times n} \,|\, |a_{i,k}^{(k)}|, |a_{k,j}^{(k)}| \le |a_{k,k}^{(k)}| \text{ for all } i, j \ge k\},$$

where $S$ is some arbitrary subset of $\mathbb{C}$ (typically $\mathbb{R}$ or $\mathbb{C}$). For instance, when performing Gaussian elimination with complete pivoting on a matrix in $\mathbf{CP}_n(S)$, no pivoting is required. Furthermore, when performing Gaussian elimination with complete pivoting on a matrix in $\mathbf{GL}_n(\mathbb{C}) \cap S^{n \times n}$, the resulting permuted matrix is in $\mathbf{CP}_n(S)$. We also stress that the role of the set $S$ is to constrain only the entries of the input matrix $A$; the sub-matrices $A^{(k)}$, $k > 1$, during Gaussian elimination need not have entries in $S$ (though if $S = \mathbb{R}$, this will of course be the case). We denote the supremum of the growth factor for a set $\mathbf{X} \subset \mathbb{C}^{n \times n}$ by $g[\mathbf{X}]$, e.g., $g[\mathbf{CP}_n(\{0, 1\})]$ is the maximum growth factor of a non-singular $n \times n$ binary matrix under complete pivoting. For all sets $\mathbf{X}$ under consideration in this work, this supremum is a maximum. In figures and tables, we use $g_n$ as shorthand for $g[\mathbf{CP}_n(\mathbb{R})]$.

1.4. **Related Work.** The maximum growth factor for partial pivoting is well understood. This quantity is known to be exactly $2^{n-1}$ for $n \times n$ complex matrices, achieved by Wilkinson's famous example matrix [36, p.212] (see [18] for all such real matrices). For complete pivoting, much less is known. A classical result, due to Wilkinson, bounds the growth factor using only Hadamard's inequality [35, Equation 4.15], and produces the estimate

$$g[\mathbf{CP}_n(\mathbb{C})] \le \sqrt{n}\big(2\,3^{1/2}\,...\,n^{1/(n-1)}\big)^{1/2} \le 2\sqrt{n}\,n^{\ln(n)/4}. \tag{1.1}$$

Minor improvements to this estimate are possible using the inexactness of Hadamard's inequality, but to date no non-trivial improvement (say, in the exponential constant) is known, even when restricted to real numbers. This estimate has historically been considered quite pessimistic; it was thought that the growth factor for real matrices under complete pivoting is at most $n$:

**Conjecture 1.1** (Folklore?[5]). $g[\mathbf{CP}_n(\mathbb{R})] \le n$, *with equality achieved only by Hadamard matrices.*

The complex analogue of this conjecture is clearly not true, as illustrated by the dimension three example [29, 31]

$$A = \begin{pmatrix} 1 & 1 & 1 \\ 1 & z & z^{-1} \\ 1 & z^{-1} & z \end{pmatrix},$$

which, when $z = \big(-1 + 2\sqrt{2}i\big)/3$, has growth factor $16/(3\sqrt{3}) \approx 3.07$. As noted by Higham, Conjecture 1.1 was one of the most famous conjectures in numerical analysis [17]. Attempting to bound or numerically compute the growth factor for small values of $n$ was a reasonably active area of research. For instance, there are a number of proofs

---

[5]See discussion in Section 1.1 regarding attribution.

that the maximum third and fourth pivots are 2.25 and 4, respectively (see Cryer [4], Tornheim [28, 29, 30, 31], Cohen [2], and Day and Peterson [5]). Tornheim also showed that the maximum fifth pivot is bounded above by $4\frac{17}{18}$ [30, 31]. Conjecture 1.1 was eventually shown to be false in dimension 13 by Gould in IEEE double precision floating point arithmetic [15], and soon after by Edelman in exact arithmetic [7]. Since these results, very little progress has been made on the asymptotic behavior of the maximum growth factor under complete pivoting or the exact values of growth for small choices of $n$.

Rook pivoting is relatively understudied compared to partial and complete pivoting, despite, in some sense, containing the best characteristics of both methods. In practice, the expected number of comparisons required should be roughly the same order of computation as partial pivoting, see [12, 24] for empirical results and theorems of this type for certain restrictive classes of random matrices. In addition, rook pivoting has a quasi-polynomial upper bound on the maximal growth factor of

$$g\big[\mathbf{RP}_n(\mathbb{C})\big] \leq \frac{3}{2}\,n^{3\ln(n)/4}, \tag{1.2}$$

as shown by Foster [12]. Similar to complete pivoting, the gap between worst-case constructions and upper bounds is quite large.

The growth factor has also been studied in a variety of other contexts. Trefethen and Schreiber studied the average growth factor over some distributions and numerical observed that for complete pivoting the growth factor appeared to exhibit a $n^{1/2}$ type behavior [34]. Higham and Higham have given numerous examples of matrices from practical applications with order $n$ growth factor [18], and recently produced a class of random matrices with growth of order $n/\log n$ [16] (both for any pivoting strategy). Sankar, Spielman, and Teng provided a smoothed analysis of growth factor without pivoting, proving that if a matrix is perturbed, it unlikely to have large growth factor [26] (in Sankar's thesis, the more complicated case of partial pivoting was also considered [25]). Recently Huang and Tikhomirov obtained new results exploring the average case analysis for partial pivoting [19]. Parker proved that, using random butterfly matrices, any non-singular matrix can be transformed into one that does not require pivoting [22]; Peca-Medlin and Trogdon further analyzed the benefits of butterfly matrices for a variety of pivoting strategies in [23]. Townsend produced bounds for the growth factor when non-optimal pivots are used [32].

1.5. **Contributions of this paper.** In this work, we prove a number of results regarding the maximum growth factor under complete pivoting, strengthen various conjectures, provide strong evidence for some results, and perform extensive numerical computations.

Through numerical search, and stability and extrapolation results, we provide improved lower bounds for the maximum growth factor:

**Theorem 1.2.** $g\big[\mathbf{CP}_n(\mathbb{R})\big] \geq 1.0045\,n$ *for all $n > 10$, and* $\limsup_n \big(g\big[\mathbf{CP}_n(\mathbb{R})\big]/n\big) \geq 3.317$.

TABLE 1. In 1991, Gould [15] presented a Table 3.1: Maximum Growth Factors Encountered. We thought it would be of interest to present the maximum growth factors we encountered over 30 years later side by side. The blue number 13.0205 is Gould's 1991 first surprising example of a matrix with $g(A) > n$. The red numbers show that it is possible to find examples even when $n = 11$ and 12. The red and magenta numbers are improvements over previously computed results. Only the bold face black numbers are known to equal $g_n$ exactly.

| $g_n$ Known Exactly | |
|---|---|
| 1 | **1** |
| 2 | **2** |
| 3 | **2.25** |
| 4 | **4** |

| Ours = same as [5, 15] | |
|---|---|
| 5 | 4.1325 |
| 6 | 5 |
| 8 | 8 |

| Ours / Gould [15] | | |
|---|---|---|
| 7 | 6.05 | 6 |
| 9 | 8.69 | 8.4305 |
| 10 | 9.96 | 9.5294 |
| **11** | **11.05** | 10.4627 |
| **12** | **12.55** | 12 |
| 13 | **13.76** | **13.0205** |
| 14 | **15.25** | 14.5949 |
| 15 | **16.92** | 16.1078 |
| 16 | **18.46** | 18.0596 |

**Conjecture:** $g_n > n$ **iff** $n > 10$

**13.0205 = 1991 surprise**

| Ours / Gould [14] | | |
|---|---|---|
| 18 | **21.25** | 20.45 |
| 20 | **24.71** | 24.25 |
| 25 | **33.67** | 32.99 |

⟵ as documented in [7]

This is the first proof that Conjecture 1.1 is false for all $n > 10$, and also the first proof that illustrates a multiplicative gap away from $n$.

In addition, we also provide asymptotic lower bounds for rook pivoting. By noting that the set of rook pivoted matrices are closed under Kronecker products, we convert finite results into lower bounds for the exponent of the growth factor, showing that rook pivoting can exhibit super-linear growth:

**Theorem 1.3.** $g\left[\mathbf{RP}_n(\mathbb{R})\right] > \frac{1}{641} n^{1.669}$ *for all* $n \in \mathbb{N}$.

Numerical search is a key ingredient in the proofs of both Theorems 1.2 and 1.3, and our numerical results also provide insights beyond the aforementioned theorems, which we briefly summarize through the following figures and tables:

- Table 1 shows improvements compared to previously known data.
- Table 2 outlines the implications of our results for low order Hadamard matrices.
- Table 3 tabulates our numerical results for every $n = 1 : 75$ and also $n = 100$.
- Figure 2 plots the numerical values from Table 3.

The reported numerical computations were performed in Julia using the modern JuMP (Julia for Mathematical Programming) [6] package. We note that, when $n = 52$ we

TABLE 2. Hadamard matrices: For decades, Hadamard matrices, interesting in their own right, seemed relevant to the growth factor problem. Gould [15] shattered that notion with his computation for $n = 16$. We observed that the notion can already be shattered partially at $n = 8$ and fully at $n = 12$:

$n = 4$: Mathematics shows $g_4 = 4$ and the optimum is Hadamard

$n = 8$: $g_8 = 8$ remains a conjecture, however one new observation is that the matrix need not be Hadamard.

$n = 12$: We are the first to report a 12x12 matrix with $g_{12} > 12$, thus showing Hadamard matrices do not maximize growth for $n = 12$.

$n = 16$: Gould reported the discovery of a 16x16 matrix with $g_{16} > 16$, therefore Hadamard matrices do not maximize growth for $n = 16$. We observed a slightly more optimal matrix.

TABLE 3. GECP Data computed by JuMP for matrices of dimensions $n = 1 : 75$ and 100 in exact arithmetic.

| $n =$ | $g \geq \downarrow$ | | $n =$ | $g \geq \downarrow$ | | $n =$ | $g \geq \downarrow$ | | $n =$ | $g \geq \downarrow$ | | $n =$ | $g \geq \downarrow$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | | 16 | 18.46 | | 31 | 45.43 | | 46 | 85.85 | | 61 | 137.55 |
| 2 | 2 | | 17 | 19.86 | | 32 | 47.74 | | 47 | 87.54 | | 62 | 141.83 |
| 3 | 9/4 | | 18 | 21.25 | | 33 | 50.36 | | 48 | 91.44 | | 63 | 144.72 |
| 4 | 4 | | 19 | 22.85 | | 34 | 52.78 | | 49 | 94.72 | | 64 | 148.05 |
| 5 | 4.13 | | 20 | 24.71 | | 35 | 54.84 | | 50 | 97.24 | | 65 | 153.98 |
| 6 | 5 | | 21 | 26.21 | | 36 | 57.66 | | 51 | 101.82 | | 66 | 157.05 |
| 7 | 6.05 | | 22 | 28.01 | | 37 | 59.91 | | 52 | 104.61 | | 67 | 162.20 |
| 8 | 8 | | 23 | 29.72 | | 38 | 63.18 | | 53 | 108.09 | | 68 | 166.89 |
| 9 | 8.69 | | 24 | 31.63 | | 39 | 64.87 | | 54 | 111.19 | | 69 | 171.33 |
| 10 | 9.96 | | 25 | 33.67 | | 40 | 67.52 | | 55 | 114.76 | | 70 | 174.45 |
| 11 | 11.05 | | 26 | 34.96 | | 41 | 70.44 | | 56 | 118.18 | | 71 | 182.98 |
| 12 | 12.55 | | 27 | 36.88 | | 42 | 73.49 | | 57 | 121.90 | | 72 | 184.91 |
| 13 | 13.76 | | 28 | 39.05 | | 43 | 77.68 | | 58 | 126.23 | | 73 | 190.57 |
| 14 | 15.25 | | 29 | 41.46 | | 44 | 79.25 | | 59 | 129.42 | | 74 | 193.28 |
| 15 | 16.92 | | 30 | 43.40 | | 45 | 82.56 | | 60 | 134.27 | | 75 | 196.79 |

$\vdots$

| 100 | 331.71 |
|---|---|

have found a matrix for which the growth factor is greater than $2n$, and at $n = 100$ the growth factor is well above $3n$. We also found a matrix for which the growth factor with rook pivoting is 640 at $n = 48$. We discuss our methodology for the computation of these results and state a pair of natural conjectures in Subsection 1.6.

We also outline our more theoretical results:

We show that the maximum growth factor over matrices with entries restricted to a subset of $\mathbb{R}$ is nearly equal to the maximum growth factor over all real matrices
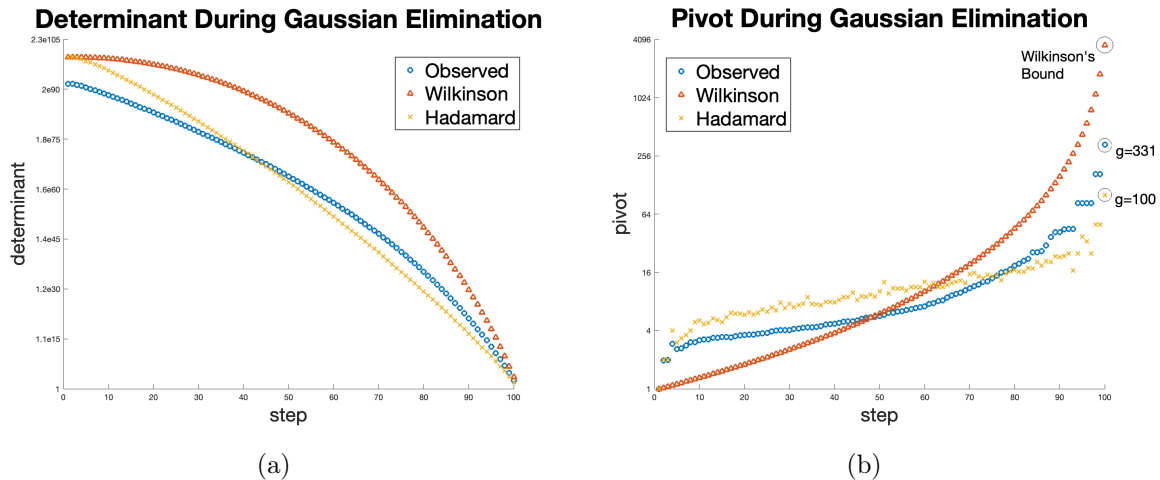
FIGURE 1. We compare the (modulus of the) determinant and pivots of $A^{(k)}$, $k = 1, ..., n$, under GECP for three examples of size $n = 100$: **Red: Wilkinson's bound**; **Yellow: a particular $n = 100$ Hadamard matrix**; **Blue: our observed maximum matrix**. (a) reveals that at least on an admittedly muted log scale, the observed determinant curve qualitatively is bending in a manner resembling Wilkinson's bound, while the Hadamard data feels qualitatively different, and thus, less relevant. (b) suggests the same conclusions as those of (a) and also suggests that "slow and steady wins the race" rather than "greedy."

**Theorem 1.4** (Simplified Version of Theorem 3.3). *For any $S \subset \mathbb{R}$, $g\big[\mathbf{CP}_{14n^2}(S)\big] \geq \big(\mathrm{diam}(S)/[2\max(S)]\big) g\big[\mathbf{CP}_n(\mathbb{R})\big]$ for all $n \in \mathbb{N}$.*

This implies that understanding the growth factor for any restricted set, say, binary matrices, is equivalent (up to polynomial factors) to understanding growth for all real matrices (i.e., if growth for binary matrices is polynomial, then it is polynomial for all matrices, and if growth for all matrices is super-polynomial, then it also is for binary matrices). We note that the $O(n^2)$ relationship is certainly pessimistic for many sets $S$ of interest; our purpose here is merely to show it is possible within small polynomial factors to find such sweeping results.

We also show that the growth factors under floating point arithmetic and exact arithmetic are nearly identical.

**Theorem 1.5** (Simplified Version of Theorem 4.2). *Let*
$$t \geq 1 + \log_\beta \big[5n^3 g^2\big[\mathbf{CP}_n(\mathbb{R})\big]\big].$$
*Then the maximum growth factor for a real $n \times n$ matrix under floating point arithmetic with base $\beta$ and mantissa length $t$ is at most $(1 + 1/n)\, g\big[\mathbf{CP}_n(\mathbb{R})\big]$.*

Theorem 4.2 treats a longstanding gap in the numerical analysis literature, a field where so much energy is devoted to the distinction between floating point and exact computations, but in the context of growth factors, this has not been analyzed to date. This

FIGURE 2. The ratio between numerically observed growth factors and matrix size for $n$ equals 1 to 75 and 100. Only the values for sizes $n = 1, 2, 3, 4$ are known mathematically to be the exact maximal growth factor though we suspect at least for the smaller values of $n$ we are achieving the maximum with our JuMP software. This data leads us to make Conjecture 1.7.

theorem provides a link between rounding error analyses that allow the unit round-off to tend to zero and those that give even more pessimistic bounds because of this difference (see, for example, [20, Theorem 2.4] and the discussion preceding it).

1.6. **Maximum Growth Factors Encountered.** With modern software and architecture we were able to find growth factors for matrices well beyond $n = 25$ as found by Gould [15] and also we were able to find larger growth for matrices as small as $n = 7$. No doubt future researchers will be able to improve our results in the same manner.

As the optimal growth problem is a constrained optimization problem it is natural to run optimization software. In 1988, Day and Peterson [5] posed the problem as a function of the $n^2$ elements of the matrix and reported some success with the FORTRAN77 nonlinear programming package NPSOL [13]. By contrast, Gould considered the advantages of posing the problem as an optimization over $n^2 + (n-1)^2 + \ldots + 1$ variables with constraints. He used the FORTRAN77 LANCELOT package that he codeveloped [3]. (LANCELOT is an acronym for "Large And Nonlinear Constrained Extended Lagrangian Optimization Techniques).

We chose to follow Gould's approach but chose to use the modern JuMP (acronym: Julia for Mathematical Programming) [6] software library with IPOPT (acronym: Interior Point Optimizer) to formulate and solve the minimization problem. The software advantage of using Julia is that the problem can be naturally formulated in a manner very similar to the mathematics [1]. The optimization engine was the COIN-OR Ipopt

(interior point optimizer) package called through Ipopt.jl. The software and results may be found in the online repository [10].

JuMP was run in parallel with 64 randomly chosen starting points on 64 separate threads and the winner, the largest growth factor, was saved. Computations were performed on a server located at MIT consisting of two AMD EPYC 7502 32-Core Hyperthreading Processors, typically using 64 of the 128 hyperthreads at a time so that others could use the machine for their own work.

We studied $n = 10$ extensively numerically and never exceeded 9.96, so we feel the evidence is very strong to state the following conjecture:

**Conjecture 1.6.** *The growth factor for complete pivoting: $g\big[\mathbf{CP}_n(\mathbb{R})\big] \geq n$ if and only if $n > 10$.*

In addition, though exact asymptotic estimates for growth factor remain elusive, we feel that we have seen sufficient numerical evidence (see Figure 2) to conjecture that the growth factor is super-linear (recall, $f(n) = \omega(g(n))$ if $\lim_{n\to\infty} f(n)/g(n) = \infty$):

**Conjecture 1.7.** $g\big[\mathbf{CP}_n(\mathbb{R})\big] = \omega(n)$.

1.7. **Remainder of Paper.** The remainder of the paper is organized as follows. In Section 2, we prove a stability lemma, which, when combined with numerical experiments and extrapolation results, imply lower bounds for both complete and rook pivoting. In Section 3, we show that the maximum growth factor for matrices with entries in an arbitrary non-trivial set $S$ is nearly as large as the maximum growth factor over all real matrices. In Section 4, we consider the growth factor in finite precision, and show that only polylogarithmically many bits (in $n$) are needed for this quantity to be at most a constant times the growth factor in exact arithmetic. In Section 5, we describe the numerical programs used to search for large growth factors, prove extrapolation results, and report our mathematically verified results. Finally, in Section 6, we study the growth factor for rook pivoting.

## 2. Key Stability Lemma: Almost Completely Pivoted is Almost Completely Pivoted

The "stability lemma" in this section is a critical technical ingredient in the majority of the theorems that follow. One immediate application follows a longstanding tradition of numerical analysis : backward error analysis. The lemma shows that if a numerical computation, such as the computations described in Sections 1 and 5, provides a computed growth factor for a "nearly" completely pivoted matrix, then there is a "nearby" matrix which has a "nearby" growth factor for complete pivoting.

For a given $\boldsymbol{\varepsilon} = (\varepsilon_1, ..., \varepsilon_{n-1}) \in \mathbb{R}^{n-1}$, $\varepsilon_i > -1$ for $i = 1, ..., n-1$, we define

$$\mathbf{CP}_n^{\boldsymbol{\varepsilon}}(S) = \{A \in \mathbf{GL}_n(\mathbb{C}) \cap S^{n \times n} \,|\, |a_{i,j}^{(k)}| \leq (1 + \varepsilon_k)|a_{k,k}^{(k)}| \text{ for all } i, j \geq k, \, (i,j) \neq (k,k)\},$$

---

**Algorithm 1** Practical Algorithm to turn optimization output to theory:
Convert $A \in \mathbf{GL}_n(\mathbb{R})$ to $B \in \mathbf{CP}_n(\mathbb{R})$

---

1: $B^{(1)} \leftarrow \mathbf{Rational}(A)$          $\triangleright$ Convert floating point matrix to rational matrix
2: **for** $k = 1 : n - 1$ **do**
3:      $B^{(k+1)} \leftarrow B^{(k)} - B_{:,k}^{(k)} B_{k,:}^{(k)} / B_{k,k}^{(k)}$
4: **end for**
5: **for** $k = (n-1) : -1 : 1$ **do**
6:      $\delta \leftarrow \max \left\{ \max_{i>k} \left( B_{i,k}^{(k)} / B_{k,k}^{(k)} \right)^2, \max_{j>k} \left( B_{k,j}^{(k)} / B_{k,k}^{(k)} \right)^2, \max_{i,j>k} \left| B_{i,j}^{(k)} / B_{k,k}^{(k)} \right| \right\}$
7:      **if** $\delta > 1$ **then**          $\triangleright$ else matrix is already CP
8:          $\delta_h \leftarrow \mathbf{Rational}\left(\mathbf{NextFloat}(\delta^{1/2})\right)$ $\triangleright$ Approx. $\delta^{1/2}$, round up, make rational
9:          $B_{:,k}^{(k)} \leftarrow \delta_h B_{:,k}^{(k)}$, $B_{k,:}^{(k)} \leftarrow \delta_h B_{k,:}^{(k)}$
10:          **if** $k > 1$ **then**
11:             **for** $i = 1 : (k-1)$ **do**
12:               $B_{k,k}^{(i)} \leftarrow B_{k,k}^{(i)} + (\delta_h^2 - 1) B_{k,k}^{(k)}$
13:               $B_{k,k+1:n}^{(i)} \leftarrow B_{k,k+1:n}^{(i)} + (\delta_h - 1) B_{k,k+1:n}^{(k)}$
14:               $B_{k+1:n,k}^{(i)} \leftarrow B_{k+1:n,k}^{(i)} + (\delta_h - 1) B_{k+1:n,k}^{(k)}$
15:             **end for**
16:          **end if**
17:      **end if**
18: **end for**
19: $B \leftarrow B^{(1)}$

---

where $S$ is some arbitrary subset of $\mathbb{C}$ (typically $\mathbb{R}$ or $\mathbb{C}$), e.g., the set of matrices that are "almost" completely pivoted (or, for $\varepsilon_k < 0$, "overly" completely pivoted) up to a multiplicative error of $\varepsilon_k$ at the $k^{th}$ step of Gaussian elimination. When $\varepsilon_1 = ... = \varepsilon_{n-1} > 0$, these sets are generally referred to as threshold-pivoted matrices.

We first prove the following lemma, showing that every matrix in $\mathbf{CP}_n^{\varepsilon}(S)$ is close to a matrix in $\mathbf{CP}_n^{\delta}(S)$, where $\varepsilon_i \geq 0 \geq \delta_i$ for all $i = 1, ..., n - 1$. The statement and proof for the more general case in which $\varepsilon_i$ and $\delta_i$ may have the same sign is similar, but is slightly more complicated and not needed for our purposes. We also give an algorithmic description of the procedure in the proof of Lemma 2.1 for the case $\boldsymbol{\delta} = 0$ in Algorithm 1, as this subroutine is a crucial part of converting numerically computed results to mathematical proofs of lower bounds.

**Lemma 2.1.** *Let* $\boldsymbol{\varepsilon} = (\varepsilon_1, ..., \varepsilon_{n-1})$ *and* $\boldsymbol{\delta} = (\delta_1, ..., \delta_{n-1})$ *satisfy* $-1 < \delta_i \leq 0 \leq \varepsilon_i$ *for* $i = 1, ..., n - 1$. *Then, for every* $A \in \mathbf{CP}_n^{\varepsilon}(S)$, *where* $S$ *equals* $\mathbb{R}$ *or* $\mathbb{C}$, *there exists a matrix* $B \in \mathbf{CP}_n^{\delta}(S)$ *such that* $b_{n,n}^{(k)} = a_{n,n}^{(k)}$ *for all* $k = 1, ..., n$, *and*

$$\left| b_{i,j}^{(k)} - a_{i,j}^{(k)} \right| \leq \max_{\min\{i,j\} \leq \ell \leq n-1} \frac{\left[ \left( \frac{1 + \varepsilon_\ell}{1 + \delta_\ell} \right)^2 - 1 \right] \left| a_{\ell,\ell}^{(\ell)} \right|}{\prod\limits_{p=\min\{i,j\}}^{\ell-1} 1 + \delta_p} + \sum_{m=\min\{i,j\}}^{\ell-1} \frac{(\varepsilon_m - \delta_m) \left| a_{m,m}^{(m)} \right|}{\prod\limits_{p=\min\{i,j\}}^{m} 1 + \delta_p}.$$

*Proof.* To construct $B \in \mathbf{CP}_n^{\delta}(S)$, we iteratively define the entries $b_{i,j}^{(k)}$, starting with $k = n$ and working backwards from $k = n$ to $k = 1$. The key to this construction is that we scale the row and column of the pivot entry of each matrix $A^{(k)}$ by a fixed multiplicative factor. This operation leaves entries $a_{i,j}^{(\ell)}$, $i, j > k$, unchanged, and so during our procedure each entry is changed at most once. The factor depends on both the maximum magnitude entry $|a_{i,j}^{(k)}|$ over all $i, j > k$ and the maximum over $i = k, j > k$ and $j = k, i > k$. This allows error to propagate additively rather than multiplicatively.

Let $B^{(n)} := A^{(n)}$ and

$$B^{(k)} := \begin{pmatrix} (1 + \gamma_k) \, a_{k,k}^{(k)} & \sqrt{1 + \gamma_k} \, A_{k,k+1:n}^{(k)} \\ \sqrt{1 + \gamma_k} \, A_{k+1:n,k}^{(k)} & A_{k+1:n,k+1:n}^{(k)} + B^{(k+1)} - A^{(k+1)} \end{pmatrix}$$

for $k = 1, ..., n-1$, where $\gamma_n := 0$ and

$$\gamma_k := \max \left\{ \left( \frac{1 + \varepsilon_k}{1 + \delta_k} \right)^2 - 1, \frac{\varepsilon_k - \delta_k + \max_{i>k} \gamma_i |a_{i,i}^{(i)}| / |a_{k,k}^{(k)}|}{1 + \delta_k} \right\}.$$

The quantity $\gamma_k |a_{k,k}^{(k)}|$ is monotonically decreasing with $k$ (as $\delta_k \leq 0$), and so we may equivalently write

$$\gamma_k |a_{k,k}^{(k)}| = \max \left\{ \left[ \left( \frac{1 + \varepsilon_k}{1 + \delta_k} \right)^2 - 1 \right] |a_{k,k}^{(k)}|, \frac{(\varepsilon_k - \delta_k)|a_{k,k}^{(k)}| + \gamma_{k+1}|a_{k+1,k+1}^{(k)}|}{1 + \delta_k} \right\},$$

or

$$\gamma_k |a_{k,k}^{(k)}| = \max_{\ell \geq k} \left[ \left( \frac{1 + \varepsilon_\ell}{1 + \delta_\ell} \right)^2 - 1 \right] |a_{\ell,\ell}^{(\ell)}| \prod_{p=k}^{\ell-1} \frac{1}{1 + \delta_p} + \sum_{m=k}^{\ell-1} (\varepsilon_m - \delta_m)|a_{m,m}^{(m)}| \prod_{p=k}^{m} \frac{1}{1 + \delta_p}.$$

Our definitions of $B^{(k)}$ are consistent with one another, as

$$B_{k+1:n,k+1:n}^{(k)} - \frac{B_{k+1:n,k}^{(k)} B_{k,k+1:n}^{(k)}}{b_{k,k}^{(k)}} = A_{k+1:n,k+1:n}^{(k)} + B^{(k+1)} - A^{(k+1)} - \frac{A_{k+1:n,k}^{(k)} A_{k,k+1:n}^{(k)}}{a_{k,k}^{(k)}} = B^{(k+1)}.$$

Furthermore, as $\varepsilon_k \leq \sqrt{1 + \gamma_k}$, for $i \leq j$ ($j > i$ is similar),

$$|b_{i,j}^{(k)} - a_{i,j}^{(k)}| = |b_{i,j}^{(i)} - a_{i,j}^{(i)}| \leq \max\{\gamma_i |a_{i,i}^{(i)}|, (\sqrt{1 + \gamma_i} - 1)|a_{i,j}^{(i)}|\}$$
$$\leq \max\{\gamma_i, (1 + \varepsilon_i)(\sqrt{1 + \gamma_i} - 1)\}|a_{i,i}^{(i)}|$$
$$= \gamma_i |a_{i,i}^{(i)}|.$$

What remains is to verify that $B \in \mathbf{CP}_n^{\delta}(S)$. We proceed by induction from $k = n - 1$ to $k = 1$. We need only consider entries in the lower right block of $B^{(k)}$, as, for $i > k$,

$$|b_{i,k}^{(k)}| = \sqrt{1 + \gamma_k}|a_{i,k}^{(k)}| \leq (1 + \varepsilon_k)\sqrt{1 + \gamma_k}|a_{k,k}^{(k)}| = \frac{1 + \varepsilon_k}{\sqrt{1 + \gamma_k}}|b_{k,k}^{(k)}| \leq (1 + \delta_k)|b_{k,k}^{(k)}|,$$

and the same bound holds for $b_{k,j}^{(k)}$, $j > k$.

When $k = n - 1$,

$$|b_{n,n}^{(n-1)}| = |a_{n,n}^{(n-1)}| \leq (1 + \varepsilon_{n-1})|a_{n-1,n-1}^{(n-1)}| = \frac{1 + \varepsilon_{n-1}}{1 + \gamma_{n-1}}|b_{n-1,n-1}^{(n-1)}| \leq (1 + \delta_{n-1})|b_{n-1,n-1}^{(n-1)}|.$$

Suppose the statement holds for $k = \ell+1, ..., n-1$, $\ell < n-1$, and consider $b_{i,j}^{(\ell)}$, $\ell < i \leq j$ ($\ell < j < i$ is similar). We have

$$
\begin{aligned}
|b_{i,j}^{(\ell)}| &\leq |a_{i,j}^{(\ell)}| + |b_{i,j}^{(i)} - a_{i,j}^{(i)}| \\
&\leq (1 + \varepsilon_\ell)|a_{\ell,\ell}^{(\ell)}| + \gamma_i|a_{i,i}^{(i)}| \\
&= |a_{\ell,\ell}^{(\ell)}|\big((1 + \varepsilon_\ell) + \gamma_i|a_{i,i}^{(i)}|/|a_{\ell,\ell}^{(\ell)}|\big) \\
&= (1 + \delta_\ell)|b_{\ell,\ell}^{(\ell)}|\frac{(1 + \varepsilon_\ell) + \gamma_i|a_{i,i}^{(i)}|/|a_{\ell,\ell}^{(\ell)}|}{(1 + \delta_\ell)(1 + \gamma_\ell)} \\
&\leq (1 + \delta_\ell)|b_{\ell,\ell}^{(\ell)}|,
\end{aligned}
$$

and therefore $B \in \mathbf{CP}_n^\delta(S)$.                                    $\square$

A simpler, but weaker version of the above result, relating the maximum growth factor under threshold complete pivoting to that of complete pivoting, is as follows.

**Corollary 2.2.** *Let $S$ equal $\mathbb{R}$ or $\mathbb{C}$, and $\boldsymbol{\varepsilon} = (\epsilon, ..., \epsilon)$, $\epsilon > 0$. Then*

$$
g\big[\mathbf{CP}_n(S)\big] \geq \frac{g\big[\mathbf{CP}_n^\varepsilon(S)\big]}{1 + \epsilon(2 + \epsilon)g\big[\mathbf{CP}_{n-1}^\varepsilon(S)\big] + \epsilon\sum_{i=1}^{n-2} g\big[\mathbf{CP}_i^\varepsilon(S)\big]}
$$

*for all $n \in \mathbb{N}$.*

*Proof.* The result follows from choosing an $A \in \mathbf{CP}_n^\varepsilon(S)$ that achieves the maximum growth factor and choosing $B \in \mathbf{CP}_n(S)$ from Lemma 2.1 as a lower bound for $g\big[\mathbf{CP}_n(S)\big]$. By Lemma 2.1, we have

$$
g(A) = \frac{|a_{n,n}^{(n)}|}{\max_{i,j}|a_{i,j}|} = \frac{|b_{1,1}|}{\max_{i,j}|a_{i,j}|}\, g(B) \leq \frac{|b_{1,1}|}{|a_{1,1}|}\, g(B) = (1 + \gamma_1)\, g(B),
$$

where

$$
\gamma_1 \leq \max_{\ell<n} \epsilon(2 + \epsilon)\frac{|a_{\ell,\ell}^{(\ell)}|}{|a_{1,1}|} + \epsilon\sum_{m=1}^{\ell-1}\frac{|a_{m,m}^{(m)}|}{|a_{1,1}|} \leq \epsilon(2 + \epsilon)g\big[\mathbf{CP}_{n-1}^\varepsilon(S)\big] + \epsilon\sum_{m=1}^{n-2} g\big[\mathbf{CP}_m^\varepsilon(S)\big].
$$

$\square$

In addition to being a crucial ingredient for our results, Corollary 2.2 also has some historical significance. This result, and the associated algorithm illustrates a way to convert almost completely pivoted matrices into matrices that are completely pivoted, without losing much in the growth factor. This has key similarities to Edelman's exact arithmetic extension of Gould's finite precision counterexample to Conjecture 1.1, and provides some answers to Edelman's perturbation question for growth factor [7].

## 3. Growth Factor for Constrained Entries

In this section, we study the maximum growth factor of matrices in $\mathbf{GL}_n(\mathbb{R}) \cap S^{n \times n}$ when $S$ is a small set (e.g., $\{0, 1\}$). In particular, we aim to show that the maximum growth factor for matrices with entries restricted to some subset $S \subset \mathbb{R}$ is nearly the same as the growth factor over $\mathbb{R}$, up to a quadratic factor in the input $n$. To do so, we proceed as follows: First, we show that the maximum growth factor for matrices at least some prescribed distance from the boundary is almost as large as the maximum growth factor over the entire set (Lemma 3.1). Combining this result with the stability lemma of the previous section (Lemma 2.1) produces a lower bound for the maximum growth factor of sets of matrices that cover $\mathbf{CP}_n(S)$ sufficiently well (Lemma 3.2). Finally, using this lower bound, we show that if our restricted set $S$ is non-trivial (i.e., $|S| > 1$), then we can almost achieve the maximum growth factor, up to a quadratic factor in $n$.

We begin by characterizing a subset of $\mathbf{CP}_n(S)$ that is stable under entry-wise perturbations of size at most $\varepsilon$, i.e., matrices $A$ such that $\{B \in S^{n \times n} \mid |a_{i,j} - b_{i,j}| \leq \varepsilon\} \subset \mathbf{CP}_n(S)$. We have the following lemma.

**Lemma 3.1.** *Let $A \in \mathbf{CP}_n(S)$, $S$ equal $\mathbb{R}$ or $\mathbb{C}$, and $\varepsilon > 0$. If $|a_{i,j}^{(k)}| \leq |a_{k,k}^{(k)}| - 2 \times 4^{k-1}\varepsilon$ for all $i, j = k, ..., n$ (except $i = j = k$), $k = 1, ..., n-1$, then*

$$\{B \in S^{n \times n} \mid |a_{i,j} - b_{i,j}| \leq \varepsilon\} \subset \mathbf{CP}_n(S),$$

*and*

$$g\big[\{B \in S^{n \times n} \mid |a_{i,j} - b_{i,j}| \leq \varepsilon\}\big] \geq g(A) - \varepsilon(4^{n-1} + g(A))/|a_{1,1}|.$$

*Proof.* Let $B \in S^{n \times n}$ satisfy $b_{i,j}^{(1)} = a_{i,j}^{(1)} + \theta_{i,j}^{(1)}$, where $|\theta_{i,j}^{(1)}| \leq \varepsilon$. Then

$$b_{i,j}^{(2)} = \left[(a_{i,j}^{(1)} + \theta_{i,j}^{(1)}) - \frac{(a_{i,1}^{(1)} + \theta_{i,1}^{(1)})(a_{1,j}^{(1)} + \theta_{1,j}^{(1)})}{(a_{1,1}^{(1)} + \theta_{1,1}^{(1)})}\right] + a_{i,j}^{(2)} - \left[a_{i,j}^{(1)} - \frac{a_{i,1}^{(1)} a_{1,j}^{(1)}}{a_{1,1}^{(1)}}\right] = a_{i,j}^{(2)} + \theta_{i,j}^{(2)},$$

where

$$\theta_{i,j}^{(2)} := \theta_{i,j}^{(1)} + \theta_{1,1}^{(1)} \frac{a_{i,1}^{(1)} a_{1,j}^{(1)}}{a_{1,1}^{(1)}(a_{1,1}^{(1)} + \theta_{1,1}^{(1)})} - \frac{\theta_{i,1}^{(1)} a_{1,j}^{(1)} + \theta_{1,j}^{(1)} a_{i,1}^{(1)} + \theta_{i,1}^{(1)} \theta_{1,j}^{(1)}}{a_{1,1}^{(1)} + \theta_{1,1}^{(1)}}.$$

Since $|a_{i,1}^{(1)}|, |a_{1,j}^{(1)}| \leq |a_{1,1}^{(1)}| - 2\varepsilon < |a_{1,1}^{(1)}| - \varepsilon \leq |a_{1,1}^{(1)} + \theta_{1,1}^{(1)}|$, we have

$$|\theta_{i,j}^{(2)}| \leq \varepsilon\left(1 + \frac{2|a_{1,j}^{(1)}| + |a_{i,1}^{(1)} + \theta_{i,1}^{(1)}|}{|a_{1,1}^{(1)} + \theta_{1,1}^{(1)}|}\right) \leq 4\varepsilon.$$

Repeating this estimate for $k = 3, ..., n$ with $\varepsilon$ replaced by $4^{k-2}\varepsilon$, we have $|a_{i,j}^{(k)} - b_{i,j}^{(k)}| \leq 4^{k-1}\varepsilon$ for all $i, j, k$. Suppose that $g(A)$ is achieved by entry $a_{i,j}^{(k)}$. Then

$$g(B) \geq \frac{|a_{i,j}^{(k)}| - 4^{k-1}\varepsilon}{|a_{1,1}| + \varepsilon} = g(A) - \frac{4^{k-1}\varepsilon|a_{1,1}| + \varepsilon|a_{i,j}^{(k)}|}{(|a_{1,1}| + \varepsilon)|a_{1,1}|} \geq g(A) - \varepsilon(4^{n-1} + g(A))/|a_{1,1}|.$$

$\square$

Combining Lemmas 2.1 and 3.1, we are now prepared to prove a lemma regarding the maximum growth factor over sets that cover $\mathbb{R}^{n \times n}$ (or $\mathbb{C}^{n \times n}$) sufficiently well.

**Lemma 3.2.** *Let $n > 1$, $0 < \varepsilon < 2^{-(2n-1)}$, and $S$ equal $\mathbb{R}$ or $\mathbb{C}$. Let $X \subset S^{n \times n}$ be a subset such that, for all $A \in \mathbf{CP}_n(S)$, there exists an $\alpha \in S$ and $B \in X$ satisfying $|a_{i,j} - \alpha\, b_{i,j}| \leq \varepsilon |a_{1,1}|$ for all $i,j = 1, ..., n$. Then $\mathbf{CP}_n(S) \cap X$ is non-empty and*

$$g\big[\mathbf{CP}_n(S) \cap X\big] \geq \big(1 - \varepsilon n 4^{n-1} g\big[\mathbf{CP}_n(S)\big]/(2\varepsilon; 4)_n\big)\, g\big[\mathbf{CP}_n(S)\big],$$

*where $(\cdot\, ; \cdot)_n$ is the q-Pochhammer symbol.*

*Proof.* The main idea of the proof is as follows. We consider a matrix $A \in \mathbf{CP}_n(S)$, $a_{1,1} = 1$, that maximizes growth factor (i.e., $g(A) = g\big[\mathbf{CP}_n(S)\big]$) and, using Lemma 2.1 applied to $\mathbf{CP}_n(S)$ and $\mathbf{CP}_n^{\boldsymbol{\delta}}(S)$ for $\boldsymbol{\delta}$ entry-wise negative, find a nearby matrix $C \in \mathbf{CP}_n^{\boldsymbol{\delta}}(S)$. Then, we find a matrix $B \in X$ nearby $C$ and, using Lemma 3.1, conclude that $B \in \mathbf{CP}_n(S)$. Finally, using the bounds on $|a_{i,j}^{(k)} - c_{i,j}^{(k)}|$ and $|b_{i,j} - c_{i,j}|$ we argue that $g(B)$ is fairly large.

So that we may apply Lemma 3.1, we define $\delta_k = -2 \times 4^{k-1}\varepsilon$ and let $C \in \mathbf{CP}_n^{\boldsymbol{\delta}}(S)$ be the matrix resulting from the proof of Lemma 2.1. Because $A$ maximizes $g(A)$, $|a_{k,k}^{(k)}| \geq 1$ for $k = 1, ..., n$ and therefore $|c_{k,k}^{(k)}| \geq 1$ for $k = 1, ..., n$ as well. In this case, $C$ satisfies

$$|c_{i,j}^{(k)}| \leq (1 + \delta_k)|c_{k,k}^{(k)}| = |c_{k,k}^{(k)}| - 2 \times 4^{k-1}\varepsilon|c_{k,k}^{(k)}| \leq |c_{k,k}^{(k)}| - 2 \times 4^{k-1}\varepsilon,$$

and so, by Lemma 3.1 combined with our lemma hypothesis, there exists a matrix $B \in \mathbf{CP}_n(S) \cap X$ (w.l.o.g. $\alpha = 1$) with $|b_{i,j} - c_{i,j}| \leq \varepsilon$.

What remains is to bound the differences $|a_{1,1} - b_{1,1}|$ and $|a_{n,n}^{(n)} - b_{n,n}^{(n)}|$, and compute a lower bound for $g(B)$. By Lemmas 2.1 and 3.1,

$$|a_{1,1} - b_{1,1}| \leq |a_{1,1} - c_{1,1}| + |b_{1,1} - c_{1,1}|$$

$$\leq g(A)\left[\frac{(1 - 2 \times 4^{n-2}\epsilon)^{-2} - 1}{\prod_{p=1}^{n-2}(1 - 2 \times 4^{p-1}\varepsilon)} + \sum_{m=1}^{n-2} \frac{2 \times 4^{m-1}\varepsilon}{\prod_{p=1}^{m}(1 - 2 \times 4^{p-1}\varepsilon)}\right] + \epsilon$$

$$= \varepsilon\left(1 + g(A)\left[\frac{2 \times 4^{n-2}}{1 - 2 \times 4^{n-2}\varepsilon}(2\varepsilon; 4)_{n-1}^{-1} + \sum_{m=1}^{n-1} 2 \times 4^{m-1}(2\varepsilon; 4)_m^{-1}\right]\right)$$

$$\leq \varepsilon\big(1 + 2\,n\,4^{n-2}g(A)/(2\varepsilon; 4)_n\big),$$

and

$$|a_{n,n}^{(n)} - b_{n,n}^{(n)}| \leq |a_{n,n}^{(n)} - c_{n,n}^{(n)}| + |b_{n,n}^{(n)} - c_{n,n}^{(n)}| \leq 4^{n-1}\varepsilon.$$

Therefore,

$$\frac{|b_{n,n}^{(n)}|}{|b_{1,1}|} \geq \frac{g(A) - 4^{n-1}\varepsilon}{1 + \varepsilon\big(1 + 2\,n\,4^{n-2}g(A)/(2\varepsilon;4)_n\big)}$$

$$= g(A) - \varepsilon\,\frac{4^{n-1} + g(A)\big(1 + 2\,n\,4^{n-2}g(A)/(2\varepsilon;4)_n\big)}{1 + \varepsilon\big(1 + 2\,n\,4^{n-2}g(A)/(2\varepsilon;4)_n\big)}$$

$$\geq g(A) - \varepsilon\,\big(4^{n-1} + g(A)\big(1 + 2\,n\,4^{n-2}g(A)/(2\varepsilon;4)_n\big)\big)$$

$$= g(A)\big(1 - \varepsilon(4^{n-1}/g(A) + 1 + 2\,n\,4^{n-2}g(A)/(2\varepsilon;4)_n)\big)$$

$$\geq g(A)\big(1 - \varepsilon n 4^{n-1} g(A)/(2\varepsilon;4)_n\big).$$

$\square$

The requirement on the cover that $X$ provides in the previous lemma is quite strong; for a non-trivial result we require $\varepsilon$ to be exponentially small in $n$. However, the set of $m \times m$ matrices $A^{(n-m+1)}$ resulting from many steps of Gaussian elimination applied to the set of $A \in S^{n \times n}$, for $S$ finite, does indeed provide an approximation to any $m \times m$ matrix with error exponentially small in $m$ (where $m$ is a sufficiently small fixed polynomial in $n$). We formalize this concept in the proof of the following theorem, which relates the maximum growth of $\mathbf{CP}_m(S)$, $S \subset \mathbb{R}$, $|S| > 1$, to that of $\mathbf{CP}_n(\mathbb{R})$.

**Theorem 3.3.** *If $S \subset \mathbb{R}$, then*

$$g\big[\mathbf{CP}_m(S)\big] \geq \frac{diam(S)}{2\max_{s \in S}|s|}g\big[\mathbf{CP}_n(\mathbb{R})\big] \quad \text{for all } m > 4n(3n+1).$$

*Proof.* The main idea of the proof is to build a matrix $B \in S^{m \times m}$, $m = n + p$, such that iterates $B^{(i)}$, $i = 1, ..., p$, are completely pivoted, $|b_{p+1,p+1}^{(p+1)}| \geq |b_{1,1}|$, and $B^{(p+1)}$ approximates an arbitrary $A \in \mathbf{CP}_n(\mathbb{R})$ exponentially well. If we can approximate an arbitrary $A$ up to error $2^{-3n}$, i.e., $|a_{i,j} - \alpha b_{i,j}^{(k+1)}| \leq 2^{-3n}|a_{1,1}|$ for some fixed $\alpha$, then, by Lemma 3.2 combined with Wilkinson's bound (Inequality 1.1) for $g\big[\mathbf{CP}_n(\mathbb{C})\big]$ (for $n > 1$),

$$g\big[\mathbf{CP}_m(S)\big] \geq (1 - 2^{-(n+1)}n^{\ln(n)/4+3/2}/(2^{1-3n};4)_n)g\big[\mathbf{CP}_n(\mathbb{R})\big] \geq \frac{1}{2}g\big[\mathbf{CP}_n(\mathbb{R})\big].$$

What remains is to construct the matrix $B$.

Given any $s_1, s_2 \in S$, $|s_1| < |s_2|$, and matrix $C \in \mathbf{CP}_{m-1}(\{0,1\})$, the matrix

$$B = \begin{pmatrix} s_2 & s_2\mathbf{1}^T \\ s_2\mathbf{1} & s_2\mathbf{1}\mathbf{1}^T + (s_1 - s_2)C \end{pmatrix}$$

is in $\mathbf{CP}_m(S)$ and satisfies $B^{(2)} = (s_1 - s_2)C$. Therefore, we may assume that $S = \{0,1\}$ at the cost of one step of Gaussian elimination and a multiplicative factor of $diam(S)/\max_{s \in S}|s|$ in the growth factor. However, we would like a matrix with entries in $\{0, 1/2, 1\}$. To do so, we note that three steps of Gaussian elimination applied to the

$(m-1) \times (m-1)$ block matrix

$$C = \begin{pmatrix} 1 & 1 & 0 & \mathbf{0}^T & 0 \\ 1 & 0 & 1 & \mathbf{0}^T & 0 \\ 0 & 1 & 1 & \mathbf{0}^T & 1 \\ \mathbf{0} & \mathbf{0} & x & E & y \end{pmatrix}$$

where $x \in \{0,1\}^{m-4}$, $E \in \{0,1\}^{(m-4)\times(m-5)}$, and $y \in \{0,1\}^{m-4}$, produces a $(m-4) \times (m-4)$ matrix with its first $m-5$ columns given by $E$ and its last column given by $y - x/2$. Performing this operation $\ell$ times produces a $\ell \times \ell$ $\{0,1/2,1\}$ matrix, where $\ell$ must be such that $4\ell + 1 \le m$. We are now prepared to approximate an arbitrary matrix $A \in \mathbf{CP}_n(\mathbb{R})$ using matrices in $\mathbf{CP}_\ell(\{0,1/2,1\})$. Suppose (w.l.o.g.) that $a_{1,1} = 1$, and let $r_{i,j,k}$ denote the $k^{th}$ bit in the binary expansion of $\text{ceil}(a_{i,j}) - a_{i,j}$ (we write $-1$ as $-0.\bar{1}$ in binary), and set $r_{i,j,0}$ to be the integer part of $a_{i,j}$ (i.e., either 0 or 1). To obtain an approximation of $A$ of order $2^{-3n}$, we set $\ell = 3n^2 + n$ and define $E$ as follows

$$E = \begin{pmatrix} I & \mathbf{0} & \cdots & \mathbf{0} & \frac{1}{2}I \\ \frac{1}{2}I & I & \ddots & \vdots & \frac{1}{2}I \\ \vdots & \ddots & \ddots & \mathbf{0} & \vdots \\ \frac{1}{2}I & \cdots & \frac{1}{2}I & I & \frac{1}{2}I \\ R_1 & R_2 & \cdots & R_{3n} & R_0 \end{pmatrix},$$

where each block is $n \times n$, and $R_k = (r_{i,j,k})_{i,j=1}^n$ for $k = 0, 1, ..., 3n$. After $n$ steps of Gaussian elimination, we have

$$E^{(n+1)} = \begin{pmatrix} I & \mathbf{0} & \cdots & \mathbf{0} & \frac{1}{4}I \\ \frac{1}{2}I & I & \ddots & \vdots & \frac{1}{4}I \\ \vdots & \ddots & \ddots & \mathbf{0} & \vdots \\ \frac{1}{2}I & \cdots & \frac{1}{2}I & I & \frac{1}{4}I \\ R_2 & R_3 & \cdots & R_{3n} & R_0 - \frac{1}{2}R_1 \end{pmatrix},$$

and finally, after $3n^2$ steps we have that $E^{(3n^2+1)} = R_0 - \frac{1}{2}R_1 - \frac{1}{4}R_2 - ... - \frac{1}{2^{3n}}R_{3n}$ and approximates $A$ up to error $2^{-3n}$. We have $\ell = 3n^2 + n$ and require $4\ell + 1 \le m$, so we set $m = 4n(3n+1) + 1$. $\square$

A similar result (with a worse multiplicative constant) holds for $\mathbb{C}$ given a set $S$ which either contains $\{0, 1, i\}$, or can be converted to such a set after relatively few iterates of Gaussian elimination (e.g., $\{-1, 1, 1+i\}$). We leave the details to the motivated reader.

## 4. Growth Factor in Floating Point Arithmetic

In this section, we aim to bound the growth factor encountered in practice in floating point arithmetic. The term "growth factor" in the literature is used ambiguously to refer to two closely related quantities: growth factor under exact arithmetic or under floating point arithmetic, leading to some confusion. The exact case is clear, and shows up in theoretical discussions. The floating point arithmetic case, by contrast refers to the largest element (in absolute value) seen during a floating point computation. As

previously mentioned in Section 1, error estimates for Gaussian elimination typically involve the growth factor under floating point arithmetic rather than exact arithmetic. In this section, we show that when using sufficiently high precision ($\omega(\log^2 n)$ bits), the maximum growth factor for exact and floating point arithmetic are identical up to a $1 + o(1)$ multiplicative factor (Theorem 4.2).

We consider the maximum growth factor when performing Gaussian elimination in base $\beta$ with $t$ digits of precision. For simplicity, we ignore issues of overflow and underflow. Here, we focus exclusively on real-valued matrices, but the analogous theorem for complex matrices follows quickly from the below analysis by simply adjusting the error due to multiplication and division for a given base and mantissa. We leave further details to the interested reader. Under floating point arithmetic, the procedure of Gaussian elimination is given by

$$\hat{a}_{i,j}^{(1)} := a_{i,j}(1 + \phi_{i,j}^{(0)}) \qquad \text{for} \quad i, j = 1, ..., n,$$

$$\hat{a}_{i,j}^{(k+1)} := \big[\hat{a}_{i,j}^{(k)} - s_{i,k}\hat{a}_{k,j}^{(k)}(1 + \theta_{i,j}^{(k)})\big](1 + \phi_{i,j}^{(k)}) \quad \text{for} \quad i, j = k, ..., n, \ k = 1, ..., n-1.$$

where

$$s_{i,k} = \frac{\hat{a}_{i,k}^{(k)}}{\hat{a}_{k,k}^{(k)}}(1 + \varphi_{i,k}),$$

and $|\theta_{i,j}^{(k)}|, |\phi_{i,j}^{(k)}|, |\varphi_{i,k}| \leq u := \beta^{1-t}/2$ for all $i, j, k$ ($u$ is commonly referred to as the unit round-off). When partial pivoting is employed, we may assume that $|s_{i,k}| \leq 1$ and $|s_{i,k}(1 + \theta_{i,j}^{(k)})| \leq 1$ for all $i, j, k$. Similar to the sets $\mathbf{CP}_n(S)$ and $\mathbf{PP}_n(S)$ defined in Section 1, we define

$$\widehat{\mathbf{CP}}_n(S) = \{A \in \mathbf{GL}_n(\mathbb{C}) \cap S^{n \times n} \,|\, \hat{a}_{k,k}^{(k)} \neq 0 \text{ for all } k, \ |\hat{a}_{i,j}^{(k)}| \leq |\hat{a}_{k,k}^{(k)}| \text{ for all } i, j \geq k\},$$

$$\widehat{\mathbf{PP}}_n(S) = \{A \in \mathbf{GL}_n(\mathbb{C}) \cap S^{n \times n} \,|\, \hat{a}_{k,k}^{(k)} \neq 0 \text{ for all } k, \ |\hat{a}_{i,k}^{(k)}| \leq |\hat{a}_{k,k}^{(k)}| \text{ for all } i \geq k\}.$$

To avoid a proliferation of indices, here and in what follows the dependence of the above sets and the growth factor on $\beta$ and $t$ is implicit. We note that, for any partially pivoted matrix, we may assume that $|s_{i,k}| \leq 1$ for all $i, k$. The growth factor under finite arithmetic is denoted by

$$G(A) := \frac{\max_{i,j,k} |\hat{a}_{i,j}^{(k)}|}{\max_{i,j} |\hat{a}_{i,j}^{(1)}|},$$

and we define $G[\mathbf{X}]$ to be the maximum growth factor under finite arithmetic (with base $\beta$ and length $t$ mantissa) over all matrices in $\mathbf{X}$. The quantity $G[\mathbf{X}]$ is a key ingredient in stability theorems of Gaussian elimination (see [20, Theorem 2.6] or [17, Theorem 9.5]). In general, the best known bounds for partial, rook, and complete pivoting is given by

$$G(A) \leq \big[1 + (1 + u)^2\big]^{n-1} = 2^{n-1} + O(nu), \tag{4.1}$$

and when $\beta = 2$, this bound can simply be replaced by $2^{n-1}$ (see [20, Section 1.2] for details). For rook and complete pivoting, $2^{n-1}$ is much more pessimistic than Inequalities 1.1 and 1.2 for exact arithmetic.

As the mantissa length $t$ tends to infinity, intuitively, the maximum growth factor under floating point arithmetic will converge to its exact arithmetic counterpart. However,

given a single matrix, the growth factor in floating point can be very different from exact arithmetic due to "near ties" causing the elimination to follow a different branch. That branch, however, is the exact branch of some nearby matrix, as the following lemma illustrates (for partial pivoting).

**Lemma 4.1.** *For every $A \in \widehat{\mathbf{PP}}_n(\mathbb{R})$, there exists a matrix $B \in \mathbf{PP}_n(\mathbb{R})$ with $b_{i,j}^{(k)} = \hat{a}_{i,j}^{(k)}$ for $i = k$ or $j = k$, and*

$$\left| \hat{a}_{i,j}^{(k)} - b_{i,j}^{(k)} \right| \le u \sum_{\ell=k}^{\min\{i,j\}-1} \left[ |\hat{a}_{i,j}^{(\ell)}| + |\hat{a}_{\ell,j}^{(\ell)}|(3+u) \right]$$

*for all $i, j = k, ..., n$ and $k = 1, ..., n-1$.*

*Proof.* The main idea is to iteratively update the lower right block of each matrix $\hat{A}^{(k)}$ so that successive matrices agree exactly, i.e., $B^{(k+1)} = B_{k+1:n}^{(k)} - B_{k+1:n,k}^{(k)} B_{k,k+1:n}^{(k)} / b_{k,k}^{(k)}$. To this end, we iteratively define $B$ so that $B^{(n)} = \hat{A}^{(n)}$ and

$$B^{(k)} = \begin{pmatrix} \hat{a}_{k,k}^{(k)} & \hat{A}_{k,:}^{(k)} \\ \hat{A}_{:,k}^{(k)} & B^{(k+1)} + \hat{A}_{k+1:n,k}^{(k)} \hat{A}_{k,k+1:n}^{(k)} / \hat{a}_{k,k}^{(k)} \end{pmatrix} \quad \text{for } k = 1, ..., n-1.$$

Clearly, successive iterates of $B$ agree with each other, and $b_{i,j}^{(k)} = \hat{a}_{i,j}^{(k)}$ for $i = k$ or $j = k$. What remains is to bound the error in the lower right block. Consider the entry $b_{i,j}^{(k)}$, where $i, j > k$ and let $m = \min\{i,j\}$. We have

$$\begin{aligned}
b_{i,j}^{(k)} &= \hat{a}_{i,j}^{(m)} + \sum_{\ell=k}^{m-1} \hat{a}_{i,\ell}^{(\ell)} \hat{a}_{\ell,j}^{(\ell)} / \hat{a}_{\ell,\ell}^{(\ell)} \\
&= \left[ \hat{a}_{i,j}^{(m-1)} - s_{i,m-1} \hat{a}_{m-1,j}^{(m-1)} (1 + \theta_{i,j}^{(m-1)}) \right] (1 + \phi_{i,j}^{(m-1)}) \\
&\quad + \hat{a}_{m-1,j}^{(m-1)} \left( s_{i,m-1} - \varphi_{i,m-1} \hat{a}_{i,m-1}^{(m-1)} / \hat{a}_{m-1,m-1}^{(m-1)} \right) + \sum_{\ell=k}^{m-2} \hat{a}_{i,\ell}^{(\ell)} \hat{a}_{\ell,j}^{(\ell)} / \hat{a}_{\ell,\ell}^{(\ell)} \\
&= \left[ \phi_{i,j}^{(m-1)} \hat{a}_{i,j}^{(m-1)} - s_{i,m-1} \hat{a}_{m-1,j}^{(m-1)} (\theta_{i,j}^{(m-1)} + \phi_{i,j}^{(m-1)} + \theta_{i,j}^{(m-1)} \phi_{i,j}^{(m-1)}) \right. \\
&\quad \left. - \hat{a}_{m-1,j}^{(m-1)} \varphi_{i,m-1} \hat{a}_{i,m-1}^{(m-1)} / \hat{a}_{m-1,m-1}^{(m-1)} \right] + \left[ \hat{a}_{i,j}^{(m-1)} + \sum_{\ell=k}^{m-2} \hat{a}_{i,\ell}^{(\ell)} \hat{a}_{\ell,j}^{(\ell)} / \hat{a}_{\ell,\ell}^{(\ell)} \right].
\end{aligned}$$

Repeating this procedure, we have

$$b_{i,j}^{(k)} = \hat{a}_{i,j}^{(k)} + \sum_{\ell=k}^{m-1} \left[ \phi_{i,j}^{(\ell)} \hat{a}_{i,j}^{(\ell)} - s_{i,\ell} \hat{a}_{\ell,j}^{(\ell)} (\theta_{i,j}^{(\ell)} + \phi_{i,j}^{(\ell)} + \theta_{i,j}^{(\ell)} \phi_{i,j}^{(\ell)}) - \hat{a}_{\ell,j}^{(\ell)} \varphi_{i,\ell} \hat{a}_{i,\ell}^{(\ell)} / \hat{a}_{\ell,\ell}^{(\ell)} \right].$$

Because our matrix is partially pivoted, $|s_{i,\ell}|$ and $|\hat{a}_{i,\ell}^{(\ell)}| / |\hat{a}_{\ell,\ell}^{(\ell)}|$ are at most one, and so

$$\left| \hat{a}_{i,j}^{(k)} - b_{i,j}^{(k)} \right| \le u \sum_{\ell=k}^{\min\{i,j\}-1} \left[ |\hat{a}_{i,j}^{(\ell)}| + |\hat{a}_{\ell,j}^{(\ell)}|(3+u) \right].$$

$\square$

By combining the above lemma with Lemma 2.1, we obtain a bound on growth factor for complete pivoting.

**Theorem 4.2.** *Let* $0 < C < 1$ *and*

$$t \geq 1 + \log_\beta \left[ \frac{(1+C)(4+5C)}{C} \sum_{m=1}^{n-1} \sum_{\ell=1}^{n-m} g\big[\mathbf{CP}_\ell(\mathbb{R})\big] g\big[\mathbf{CP}_m(\mathbb{R})\big] \right].$$

*Then* $G\big[\widehat{\mathbf{CP}}_n(\mathbb{R})\big] \leq (1+C)\, g\big[\mathbf{CP}_n(\mathbb{R})\big].$

*Proof.* Suppose $A \in \widehat{\mathbf{CP}}_n(\mathbb{R})$ maximizes growth, i.e., $G(A) = G\big[\widehat{\mathbf{CP}}_n(\mathbb{R})\big]$, and let $B \in \mathbf{PP}_n(\mathbb{R})$ be a matrix satisfying $b_{k,k}^{(k)} = \hat{a}_{k,k}^{(k)}$ for all $k$, and the bounds of Lemma 4.1. Then $B \in \mathbf{CP}_n^{\boldsymbol{\varepsilon}}(\mathbb{R})$, $\boldsymbol{\varepsilon} = (\varepsilon_1, ..., \varepsilon_{n-1})$, for $\epsilon_k := (4+u)u \sum_{\ell=1}^{n-k} G\big[\widehat{\mathbf{CP}}_\ell(\mathbb{R})\big]$, as

$$|b_{i,j}^{(k)}| \leq |\hat{a}_{i,j}^{(k)}| + u \sum_{\ell=k}^{n-1} \left[ |\hat{a}_{i,j}^{(\ell)}| + |\hat{a}_{\ell,j}^{(\ell)}|(3+u) \right] \leq |b_{k,k}^{(k)}|\left( 1 + u \sum_{\ell=k}^{n-1} \frac{|\hat{a}_{i,j}^{(\ell)}|}{|\hat{a}_{k,k}^{(k)}|} + \frac{|\hat{a}_{\ell,j}^{(\ell)}|}{|\hat{a}_{k,k}^{(k)}|}(3+u) \right).$$

In addition, $G\big[\widehat{\mathbf{CP}}_n(\mathbb{R})\big] = G(A) = |b_{n,n}^{(n)}|/|b_{1,1}|$. Using Lemma 2.1 applied to $B$, we can find a matrix $C \in g\big[\mathbf{CP}_n(\mathbb{R})\big]$ that satisfies $b_{n,n}^{(n)} = c_{n,n}^{(n)}$ and

$$|c_{1,1}| = \left( 1 + \max_{\ell \leq n-1} \varepsilon_\ell(2+\varepsilon_\ell)|\hat{a}_{\ell,\ell}^{(\ell)}|/|\hat{a}_{1,1}| + \sum_{m=1}^{\ell-1} \varepsilon_m|\hat{a}_{m,m}^{(m)}|/|\hat{a}_{1,1}| \right).$$

For the sake of space, we define $\gamma := (4+u)u$, $g(n) := g\big[\mathbf{CP}_n(\mathbb{R})\big]$, and $G(n) := G\big[\widehat{\mathbf{CP}}_n(\mathbb{R})\big]$, and note that

$$G(n) \leq g(n)\left( 1 + \max_{\ell \leq n-1} \varepsilon_\ell(2+\varepsilon_\ell)G(\ell) + \sum_{m=1}^{\ell-1} \varepsilon_m G(m) \right)$$

$$\leq g(n)\left( 1 + \gamma \max_{\ell \leq n-1} G(\ell)\left[ 1 + \gamma \sum_{p=1}^{n-\ell} G(p) \right] \sum_{p=1}^{n-\ell} G(p) + \sum_{m=1}^{\ell} \sum_{p=1}^{n-m} G(p)G(m) \right) \quad (4.2)$$

$$\leq g(n)\left( 1 + \gamma \left[ 2 + \gamma \sum_{\ell=1}^{n-1} G(\ell) \right] \sum_{m=1}^{n-1} \sum_{p=1}^{n-m} G(p)G(m) \right).$$

The result follows from noting that if

$$\frac{1}{\gamma} \geq \frac{(1+C)^2}{C}\big(2 + C/2(1+C)\big) \sum_{m=1}^{n-1} \sum_{\ell=1}^{n-m} g(\ell)g(m) \quad (4.3)$$

for some $C > 0$, then $G(k) \leq (1+C)g(k)$ for all $k = 1, ..., n$. Indeed, we have $G(1) = g(1)$, and, assuming $G(\ell) \leq (1+C)g(\ell)$ for $\ell = 1, ..., k$,

$$\frac{G(k+1)}{g(k+1)} \leq 1 + \gamma(1+C)^2 \left[ 2 + \gamma(1+C) \sum_{\ell=1}^{n-1} g(\ell) \right] \sum_{m=1}^{n-1} \sum_{p=1}^{n-m} g(p)g(m)$$

$$\leq 1 + \gamma(1+C)^2 \left[ 2 + \frac{C}{2(1+C)} \right] \sum_{m=1}^{n-1} \sum_{p=1}^{n-m} g(p)g(m) \leq 1 + C.$$

TABLE 4. Lower bounds for the largest $n$ satisfying the conditions of Theorem 4.2 for IEEE 754 Double and Quadruple precision, with $C = 1/2$ (i.e., $G\big[\widehat{\mathbf{CP}}_n(\mathbb{R})\big] \leq (3/2)\,g\big[\mathbf{CP}_n(\mathbb{R})\big]$) under various assumptions on upper bounds for $g\big[\mathbf{CP}_n(\mathbb{R})\big]$.

| $g\big[\mathbf{CP}_n(\mathbb{R})\big] \leq ...$ | Double Prec. $(t = 52)$ | Quadruple Prec. $(t = 112)$ |
|---|---|---|
| $3n$ | 4188 | 137266926 |
| $n^2/2$ | 660 | 676504 |
| Inequality 1.1 | 554 | 29563 |

$\square$

The above theorem is incredibly pessimistic, but nevertheless still provides useful some useful information. First, by using Wilkinson's bound, we note that $G\big[\widehat{\mathbf{CP}}_n(\mathbb{R})\big] \leq (1 + 1/\mathrm{Poly(n)})\,g\big[\mathbf{CP}_n(\mathbb{R})\big]$ for $t = \omega(\log_\beta^2(n))$, and, under the assumption that $g\big[\mathbf{CP}_n(\mathbb{R})\big]$ is bounded by a polynomial, only a $t = \omega(\log_\beta(n))$ length mantissa is required. In Table 4, we include a number of possible bounds on the growth factor (including Wilkinson's Inequality 1.1), and list lower bounds on the largest value of $n$ for which Theorem 4.2 guarantees that $G\big[\widehat{\mathbf{CP}}_n(\mathbb{R})\big] \leq (3/2)\,g\big[\mathbf{CP}_n(\mathbb{R})\big]$.

## 5. Computer-Assisted Lower Bounds

In this section, we detail lower bounds for growth factor found using computer search, and discuss how such computer-generated matrices in finite arithmetic lead to mathematically provable lower bounds for growth factor in exact arithmetic (Theorem 5.2).

5.1. **Computer-Assisted Lower Bounds for Small Dimension.** We are indebted to the early pioneering numerical optimization given by Day & Peterson [5] and Gould [15]. We are the beneficiary of more readily usable quality software (JuMP [6]), the ready availability of faster processors, and also modern parallel computing.

Our methodology is to run 64 threads each with a random $n \times n$ starting matrix of standard normals which has rows and columns permuted so that the matrix is completely pivoted. We then normalize by dividing by the $(1, 1)$ element. Our optimization is over the $1 + 2^2 + \ldots + n^2$ elements that are seen by Gaussian Elimination as suggested by Gould [15]. Therefore the starting point requires using all of these $\approx n^3/3$ elements. We store the variables in a 3-d array $x_{i,j,k}, 1 \leq k \leq n, k \leq i, j \leq n$. Thus $k = 1$ is the original matrix, $k = 2$ is the $n - 1 \times n - 1$ matrix obtained after one step of Gaussian elimination.

We include a listing of the high level function from the online repository [10] that performs the optimization in Figure 3; it is quite easy to read as it is similar to the mathematics. We invite readers to note the six lines of code that indicate the nonlinear

```
# maximize pivot growth for real matrix with complete pivoting
function run_model(n)
    model = Model(Ipopt.Optimizer)
    indices = [ (i, j, k) for k = 1:n for i = k:n for j = k:n]

    startmatrix = randn(n, n)
    thestart = genp(geperm(startmatrix))
    thestart = genp(thestart[:,:,1] * Diagonal(sign.(thestart[k,k,k] for k = 1:n)))
    thestart ./= thestart[1,1,1]

    @variable(model, x[i=indices], start = thestart[i[1],i[2],i[3]]  ) # random starts
    for k in 1:(n - 1), i in (k + 1):n, j in (k + 1:n)
        @NLconstraint(model,
        x[(k, k, k)] *( x[(i, j, k + 1)] - x[(i, j, k)] ) +  x[(i, k, k)] * x[(k, j, k)]  == 0 )
    end
    for k = 1:n
        @constraint(model, x[(k, k, k)] ≥ 0)
    end
    @constraint(model,    x[(1, 1, 1)] == 1)
    for i in 1:n, j in 1:n
        @constraint(model, -1 ≤ x[(i, j, 1)] ≤ 1)
    end
    for k in 2:n - 1, i in k:n, j in k:n
        @constraint(model,   x[(i, j, k)] ≤ x[(k, k, k)])
        @constraint(model, -x[(k, k, k)] ≤ x[(i, j, k)] )
    end
    w = n
    @objective(model, Max, x[(w, w, w)])
    set_optimizer_attribute(model, "max_iter", 500)
    set_silent(model)
    optimize!(model)
    A = reshape((value.(x)).data[1:n^2], n, n)
    B = convert_to_cp(Rational{BigInt}.(A))
    B = B/B[1,1]
    val = genp(B)[n,n,n]
    # return the optimization and the argmax
    val,B,MOI.FEASIBLE_POINT
end
```

FIGURE 3. Our run_model function performs growth optimization from a random start. We encourage the reader to examine the constraints: they correspond to the mathematical constraints a completely pivoted matrix satisfies and are easy to read. The variable $x$ is a three-dimensional array that stores the Gaussian elimination "pyramid," e.g., $x[(i, j, k)]$ is the $(i, j)^{th}$ entry of the $k^{th}$ step of Gaussian elimination and the $(k, k, k)$ entry is the $k^{th}$ pivot.

constraints (@NLconstraint) and linear constraints (@constraint), the first one of which is the constraint of Gaussian elimination:

Very importantly we also wish to discuss the line towards the bottom that begins B = convert_to_cp(Rational...  as this line turns a floating point answer to a rigorous mathematical answer. A simple observation is that an output of optimization software is not yet a theoretical lower bound because of floating point effects. In particular it is possible that an output of a program is not completely pivoted. The examples from Gould [15, 14] were close in the floating point sense to being optimums and some minor tweaking was needed [7] for the purpose of exact mathematics. In [7] the first author asked if there would always be a nearby floating point matrix. In this paper, we show that Lemma 2.1 theoretically states there would be nearby matrix, and Algorithm 1 as embodied in the convert_to_cp function working on a rational
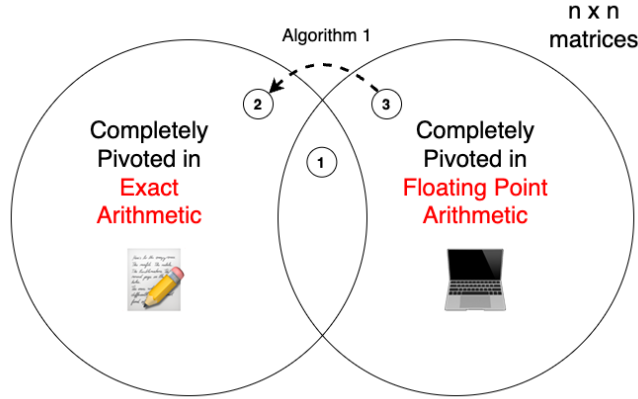
FIGURE 4. The above figure shows that we can go from matrices that are completely pivoted in floating point to matrices that are completely pivoted in exact arithmetic. Lemma 2.1 proves that this is possible and Algorithm 1 provides a pseudocode implementation (a Julia implementation may be found in the online repository [10]). For instance, Algorithm 1 has fully automated Edelman's exact arithmetic extension of Gould's finite precision counterexample to Conjecture 1.1, and provides some answers to Edelman's perturbation question for growth factor [7].

form allows us to state that our computer assisted solutions constitute exact rigorous mathematics rather than a floating point approximation.

For the smaller values of $n$, we tend to believe that the the lower bounds found may well be close to $g\big[\mathbf{CP}_n(\mathbb{R})\big]$ as we have on occasion rerun these values, and found the same answers. For larger values of $n$, we imagine that the lower bounds are just that, lower bounds.

5.2. **Provable Lower Bounds.** Next, we prove global lower bounds for the growth factor under complete pivoting by converting estimates for growth factor for small $n$ into estimates for all $n$. We begin with the following extrapolation lemma.

**Lemma 5.1.** *Let $S$ equal $\mathbb{R}$ or $\mathbb{C}$. Then*

*(i)* $g\big[\mathbf{CP}_n(S)\big]$ *is non-decreasing,*
*(ii)* $g\big[\mathbf{CP}_{2n}(S)\big] \geq 2\,g\big[\mathbf{CP}_n(S)\big]$ *for all $n \in \mathbb{N}$,*
*(iii) if* $g\big[\mathbf{CP}_n(S)\big] \geq Cn$ *for* $n = k, ..., 2k-1$, *then* $g\big[\mathbf{CP}_n(S)\big] \geq \frac{(1/k;1/2)_\infty}{1-1/k}\, Cn$ *for all* $n \geq k$, *where* $(\cdot\,;\cdot)_\infty$ *is the q-Pochhammer symbol.*

*Proof.* Properties $(i)$ and $(ii)$ follow simply from the operations

$$\begin{pmatrix} 1 & 0_n^T \\ 0_n & A \end{pmatrix} \quad \text{and} \quad A \otimes H_1, \text{ where } H_1 := \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix},$$

applied to a matrix $A \in \mathbf{CP}_n(S)$, respectively (Property $(ii)$ is also proved in [31]). If $g\big[\mathbf{CP}_n(S)\big] \geq C\,n$ for all $n \in [k, 2k)$, then by Properties $(i)$ and $(ii)$,

$$g\big[\mathbf{CP}_{2n+1}(S)\big] \geq g\big[\mathbf{CP}_{2n}(S)\big] \geq 2Cn = \frac{2n}{2n+1}\,C(2n+1) \geq \frac{2k}{2k+1}\,C(2n+1)$$

for all $n \in [k, 2k)$, i.e., $g\big[\mathbf{CP}_n(S)\big] \geq \frac{2k}{2k+1}\,Cn$ for all $n \in [k, 4k)$. Repeating this argument, we obtain the lower bound

$$g\big[\mathbf{CP}_n(S)\big] \geq Cn \prod_{i=1}^{j} \frac{2^i k}{2^i k + 1} \geq Cn \prod_{i=1}^{j} \left(1 - \frac{1}{2^i k}\right) = \frac{(1/k; 1/2)_{j+1}}{1 - 1/k}\,Cn$$

for $n \in \big[k, 2^{j+1}k\big)$, where $(\cdot\,;\cdot)_j$ is the $q$-Pochhammer symbol. Noting that $(\cdot\,;\cdot)_j$ is monotonically non-increasing with respect to $j$ for non-negative inputs of magnitude at most one completes the proof of Property $(iii)$. $\qquad\square$

Combining Lemma 5.1 with the computer-assisted (and mathematically provable) lower bounds of Table 3 immediately implies a lower bound for all values of $n$.

**Theorem 5.2** (Restatement of Theorem 1.2). $g\big[\mathbf{CP}_n(\mathbb{R})\big] \geq 1.0045\,n$ *for all* $n > 10$, *and* $\limsup_n \big(g\big[\mathbf{CP}_n(\mathbb{R})\big]/n\big) \geq 2.525$.

*Proof.* The lower bound for all $n > 10$ follows from checking $n = 11, 12, 13$ by hand and applying Property $(iii)$ of Lemma 5.1 to $k = 14$ (with $C = 1.08$). The asymptotic bound follows directly from our lower bound for $n = 100$ combined with Property $(ii)$ of Lemma 5.1. $\qquad\square$

## 6. Rook Pivoting

The majority of this work focuses on complete pivoting, due to its theoretical and practical importance. Rook pivoting by comparison is relatively understudied, yet the quasi-polynomial bound on growth factor combined with a reduced computational complexity compared to complete pivoting in practice makes this an attractive technique. Many of the results of this paper also apply to rook pivoting, sometimes leading to even stronger results. These details are left to the interested reader.

Through a stability lemma, tensor argument, and numerically computed lower bounds for a fixed value of $n$, we provide improved lower bounds for the maximum growth factor with rook pivoting. Let

$$\mathbf{RP}_n^{\varepsilon}(S) = \{A \in \mathbf{GL}_n(\mathbb{C}) \cap S^{n \times n} \,|\, |a_{i,k}^{(k)}|, |a_{k,j}^{(k)}| \leq (1 + \varepsilon_k)|a_{k,k}^{(k)}| \text{ for all } i, j \geq k\}.$$

We have the following proposition (in the spirit of Lemma 2.1).

**Proposition 6.1.** *For every* $A \in \mathbf{RP}_n^{\varepsilon}(S)$, *where* $S$ *equals* $\mathbb{R}$ *or* $\mathbb{C}$ *and* $\varepsilon_i \geq 0$ *for* $i = 1, ..., n-1$, *there exists a matrix* $B \in \mathbf{RP}_n(S)$ *such that*

$$a_{n,n}^{(k)} = b_{n,n}^{(k)} \qquad \text{for} \quad k = 1, ..., n,$$

*and*

$$\left|a_{i,j}^{(k)} - b_{i,j}^{(k)}\right| \le (2 + \varepsilon_\ell)\varepsilon_\ell \left|a_{\ell,\ell}^{(\ell)}\right|, \quad \ell := \min\{i, j\},$$

*for all* $i, j = k, ..., n, \ k = 1, ..., n - 1$.

*Proof.* Given $A \in \mathbf{RP}_n^{\varepsilon}(S)$, the result follows immediately from the construction $B^{(n)} := A^{(n)}$ and

$$B^{(k)} := \begin{pmatrix} (1 + \varepsilon_k)^2 \, a_{k,k}^{(k)} & (1 + \varepsilon_k) \, A_{k,k+1:n}^{(k)} \\ (1 + \varepsilon_k) \, A_{k+1:n,k}^{(k)} & A_{k+1:n,k+1:n}^{(k)} + B^{(k+1)} - A^{(k+1)} \end{pmatrix}$$

for $k = 1, ..., n - 1$. $\qquad\qquad\square$

Similar to Lemma 2.1, the construction of $B \in \mathbf{RP}_n(S)$ is algorithmic in nature, and this procedure (a variant of Algorithm 1) converts inexact numerically computed instances of large growth into provable lower bounds. In particular, through the combination of numerical computation and an algorithmic implementation of the procedure of Lemma 6.1, we have the following lower bound (see Subsection 5.1 and our repository [10] ).

**Proposition 6.2.** $g\big[\mathbf{RP}_{48}(\mathbb{R})\big] > 640.4861$.

Next, we prove the following extrapolation lemma, from which lower bounds for rook pivoting immediately follow.

**Lemma 6.3.** *Let* $S$ *equal* $\mathbb{R}$ *or* $\mathbb{C}$. *Then*

(i) $g\big[\mathbf{RP}_n(S)\big]$ *is non-decreasing,*
(ii) $g\big[\mathbf{RP}_{mn}(S)\big] \ge g\big[\mathbf{RP}_m(S)\big] \times g\big[\mathbf{RP}_n(S)\big]$ *for all* $m, n \in \mathbb{N}$,
(iii) *if* $g\big[\mathbf{RP}_k(S)\big] \ge k^\alpha$ *for some* $k$, *then* $g\big[\mathbf{RP}_n(S)\big] \ge k^{-\alpha}n^\alpha$ *for all* $n \in \mathbb{N}$.

*Proof.* Property $(i)$ follows from the construction $\begin{pmatrix} 1 & 0_n^T \\ 0_n & A \end{pmatrix}$. Property $(ii)$ follows from the fact that if $A \in \mathbf{RP}_m(S)$, $B \in \mathbf{RP}_n(S)$, and $S$ is closed under addition and multiplication, then $A \otimes B \in \mathbf{RP}_{mn}(S)$, where $\otimes$ is the matrix Kronecker product, which we now prove. Let $C = A \otimes B$, and, for the sake of space, define the following three auxilary matrices, consisting of $B^{(k)}$ for some $k = 2, ..., n$ and some zeros:

$$B_r^{(k)} = \begin{pmatrix} 0_{n-k+1,k-1} & B^{(k)} \end{pmatrix}, \quad B_c^{(k)} = \begin{pmatrix} 0_{k-1,n-k+1} \\ B^{(k)} \end{pmatrix}, \quad B_f^{(k)} = \begin{pmatrix} 0_{k-1,k-1} & 0_{k-1,n-k+1} \\ 0_{n-k+1,k-1} & B^{(k)} \end{pmatrix},$$

so that $B_r^{(k)} \in S^{n \times (n-k+1)}$, $B_c^{(k)} \in S^{(n-k+1) \times n}$, and $B_f^{(k)} \in S^{n \times n}$. It suffices to complete $n$ steps of Gaussian elimination, show that at each step the rook pivoting condition holds ($|c_{k,k}^{(k)}| \ge |c_{i,k}^{(k)}|, |c_{k,j}^{(k)}|$ for $k = 1, ..., n$), and note that $C^{(n+1)} = A^{(2)} \otimes B$. Initially, we have

$$C^{(1)} = A \otimes B = \begin{pmatrix} a_{1,1}B & \cdots & a_{1,m}B \\ \vdots & \ddots & \vdots \\ a_{m,1}B & \cdots & a_{m,m}B \end{pmatrix},$$

and the rook pivoting condition holds initially for any Kronecker product $A \otimes B$ of rook pivoted matrices $A$ and $B$, as

$$|a_{1,1} b_{1,1}| = |a_{1,1}|\,|b_{1,1}| \geq \max_{i,j=1,\ldots,m} \{|a_{i,1}|, |a_{1,j}|\}\ \max_{i,j=1,\ldots,n} \{|b_{i,1}|, |b_{1,j}|\}.$$

On the $k^{th}$ step of Gaussian elimination, we have

$$C^{(k)} = \begin{pmatrix} a_{1,1}B^{(k)} & a_{1,2}B_r^{(k)} & \cdots & a_{1,m}B_r^{(k)} \\ a_{2,1}B_c^{(k)} & a_{2,2}^{(2)}B + (a_{2,2} - a_{2,2}^{(2)})B_f^{(k)} & \cdots & a_{2,m}^{(2)}B + (a_{2,m} - a_{2,m}^{(2)})B_f^{(k)} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1}B_c^{(k)} & a_{m,2}^{(2)}B + (a_{m,2} - a_{m,2}^{(2)})B_f^{(k)} & \cdots & a_{m,m}^{(2)}B + (a_{m,m} - a_{m,m}^{(2)})B_f^{(k)} \end{pmatrix},$$

and still the rook pivoting condition holds, as both $A$ and $B^{(k)}$ are rook pivoted. Finally, after the $n^{th}$ step, we note that the remainder term $(a_{i,j} - a_{i,j}^{(2)})B_f^{(n)}$ disappears, as

$$(a_{i,j} - a_{i,j}^{(2)})B_f^{(n)} - \frac{a_{i,1}a_{1,j}}{a_{1,1}} \frac{B_c^{(n)}\,B_r^{(n)}}{b_{n,n}^{(n)}} = 0_{n \times n},$$

and so $C^{(n+1)} = A^{(2)} \otimes B$.

Property $(iii)$ follows quickly from Properties $(i)$ and $(ii)$. Let $n > k$ (if $n \leq k$, the result trivially holds), and let $\ell \in \mathbb{N}$ be the largest number such that $k^\ell \leq n$. We have

$$g\big[\mathbf{RP}_n(S)\big] \geq g\big[\mathbf{RP}_{k^\ell}(S)\big] \geq k^{\alpha\ell} = \big[k^\ell/n\big]^\alpha n^\alpha \geq k^{-\alpha}n^\alpha.$$

$\square$

Using Proposition 6.2 and Lemma 6.3, we obtain our desired lower bound.

**Theorem 6.4** (Restatement of Theorem 1.3). $g\big[\mathbf{RP}_n(\mathbb{R})\big] > \frac{1}{641}n^{1.669}$ *for all* $n \in \mathbb{N}$.

## Acknowledgements

## References

[1] Jeff Bezanson, Alan Edelman, Stefan Karpinski, and Viral B Shah. Julia: A fresh approach to numerical computing. *SIAM review*, 59(1):65–98, 2017.

[2] AM Cohen. A note on pivot size in Gaussian elimination. *Linear Algebra and its Applications*, 8(4):361–368, 1974.

[3] A.R. Conn, Nick Gould, and Ph.L. Toint. *LANCELOT: a Fortran package for large-scale nonlinear optimization*. Springer Series in Computational Mathematics, vol. 17. Springer Verlag (Heidelberg, New York), 1992.

[4] Colin W Cryer. Pivot size in Gaussian elimination. *Numerische Mathematik*, 12(4):335–345, 1968.

[5] Jane Day and Brian Peterson. Growth in Gaussian elimination. *The American mathematical monthly*, 95(6):489–513, 1988.

[6] Iain Dunning, Joey Huchette, and Miles Lubin. JuMP: A modeling language for mathematical optimization. *SIAM Review*, 59(2):295–320, 2017.

[7] Alan Edelman. The complete pivoting conjecture for Gaussian elimination is false. *The Mathematica Journal*, 2(2):58–61, 1992.

[8] Alan Edelman and David Friedman. A counterexample to a Hadamard matrix pivot conjecture. *Linear and Multilinear Algebra*, 44(1):53–56, 1998.

[9] Alan Edelman and Walters Mascarenhas. On the complete pivoting conjecture for a Hadamard matrix of order 12. *Linear and Multilinear Algebra*, 38(3):181–187, 1995.

[10] Alan Edelman and John Urschel. GitHub repository associated with paper. `https://github.com/alanedelman/CompletePivotingGrowth`.

[11] Leslie V. Foster. Gaussian elimination with partial pivoting can fail in practice. *SIMAX*, 15(4):1354–1362, 1994.

[12] Leslie V Foster. The growth factor and efficiency of Gaussian elimination with rook pivoting. *Journal of Computational and Applied Mathematics*, 86(1):177–194, 1997.

[13] P. E. Gill, W. Murray, M. A. Saunders, and M. H. Wright. User's guide for SOL/NPSOL: A Fortran package for nonlinear programming. Technical Report Tech. Report SOL 83-12, Stanford University Systems Optimization Laboratory, Stanford, CA, 1983.

[14] Nick Gould. personal communication, 1991.

[15] Nick Gould. On growth in Gaussian elimination with complete pivoting. *SIAM Journal on Matrix Analysis and Applications*, 12(2):354–361, 1991.

[16] Desmond J Higham, Nicholas J Higham, and Srikara Pranesh. Random matrices generating large growth in LU factorization with pivoting. *SIAM Journal on Matrix Analysis and Applications*, 42(1):185–201, 2021.

[17] Nicholas J Higham. *Accuracy and stability of numerical algorithms*. SIAM, 2002.

[18] Nicholas J Higham and Desmond J Higham. Large growth factors in Gaussian elimination with pivoting. *SIAM Journal on Matrix Analysis and Applications*, 10(2):155–164, 1989.

[19] Han Huang and Konstantin Tikhomirov. Average-case analysis of the Gaussian elimination with partial pivoting. *arXiv preprint arXiv:2206.01726*, 2022.

[20] Eugene Isaacson and Herbert Bishop Keller. *Analysis of numerical methods*. Courier Corporation, 2012.

[21] Christos Kravvaritis and Marilena Mitrouli. On the complete pivoting conjecture for Hadamard matrices: further progress and a good pivots property. *Numerical Algorithms*, 62:571–582, 2013.

[22] Douglass Stott Parker. *Random butterfly transformations with applications in computational linear algebra*. UCLA Computer Science Department, 1995.

[23] John Peca-Medlin and Thomas Trogdon. Growth factors of random butterfly matrices and the stability of avoiding pivoting. *arXiv preprint arXiv:2203.15921*, 2022.

[24] George Poole and Larry Neal. The rook's pivoting strategy. *Journal of Computational and Applied Mathematics*, 123(1-2):353–369, 2000.

[25] Arvind Sankar. *Smoothed analysis of Gaussian elimination*. PhD thesis, Massachusetts Institute of Technology, 2004.

[26] Arvind Sankar, Daniel A Spielman, and Shang-Hua Teng. Smoothed analysis of the condition numbers and growth factors of matrices. *SIAM Journal on Matrix Analysis and Applications*, 28(2):446–476, 2006.

[27] Jennifer Seberry. Google scholar search. `https://scholar.google.com/scholar?hl=en&as_sdt=0%2C47&q=jennifer+seberry+pivoting&btnG`. Accessed: 2023-03-05.

[28] Leonard Tornheim. Pivot size in Gauss reduction. *Tech Paper, Chevron Research Co., Richmond CA*, 1964.

[29] Leonard Tornheim. Maximum third pivot for Gaussian reduction. In *Tech. Report*. Calif. Res. Corp Richmond, Calif, 1965.

[30] Leonard Tornheim. A bound for the fifth pivot in Gaussian elimination. *Tech Report, Chevron Research Co., Richmond CA*, 1969.
[31] Leonard Tornheim. Maximum pivot size in Gaussian elimination with complete pivoting. *Tech Report, Chevron Research Co., Richmond CA*, 10, 1970.
[32] Alex Townsend. Gaussian elimination corrects pivoting mistakes. *arXiv preprint arXiv:1602.06602*, 2016.
[33] Alex Townsend and Lloyd N Trefethen. Gaussian elimination as an iterative algorithm. *SIAM News*, 46, 2013.
[34] Lloyd N Trefethen and Robert S Schreiber. Average-case stability of Gaussian elimination. *SIAM Journal on Matrix Analysis and Applications*, 11(3):335–360, 1990.
[35] James Hardy Wilkinson. Error analysis of direct methods of matrix inversion. *Journal of the ACM (JACM)*, 8(3):281–330, 1961.
[36] J.H. Wilkinson. *The Algebraic Eigenvalue Problem*. Clarendon Press, Oxford, 1965.

DEPARTMENT OF MATHEMATICS, MASSACHUSETTS INSTITUTE OF TECHNOLOGY, CAMBRIDGE, MA, 02139 USA.

*Email address*: edelman@mit.edu

DEPARTMENT OF MATHEMATICS, MASSACHUSETTS INSTITUTE OF TECHNOLOGY, CAMBRIDGE, MA, 02139 USA.

*Email address*: urschel@mit.edu